

# Winning Space Race with Data Science

Abdeen Mohialden  
10/2/2024



# Outline

---

- Executive Summary
- Introduction
- Methodology
- Results
- Conclusion
- Appendix

# Executive Summary

---

## 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
2. Geospatial analytics
3. 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**

Section 1

# Methodology

# Methodology

---

## 1. Data Collection

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

## 2. Data Wrangling

- Using the `.fillna()` method to remove NaN values
- Replace NaN values with mean of columns
- 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

### 1. Interactive Visual Analytics

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

### 2. 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 best parameters using GridSearchCV
- Plotting confusion matrices for each classification model
- Assessing the accuracy of each classification model

# Data Collection

---

1

- Make a GET response to the SpaceX REST API
- Convert the response to a .json file then to a Pandas DataFrame

2

- clean the data
- Define lists for data to be stored in
- retrieve data and fill the lists
- Use these lists as values in a dictionary and construct the dataset

3

- Create a Pandas DataFrame from the constructed dictionary dataset

4

- Filter the DataFrame to only include Falcon 9 launches
- Reset the FlightNumber column
- Replace missing values of PayloadMass with the mean PayloadMass value

# 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.

1

- Make a GET response to the SpaceX REST API
- Convert the response to a .json file then to a Pandas DataFrame

2

- clean the data (see 2)
- Define lists for data to be stored in
- Call custom functions (see 4) to retrieve data and fill the lists
- Use these lists as values in a dictionary and construct the dataset

3

- Create a Pandas DataFrame from the constructed dictionary dataset

4

- Filter the DataFrame to only include Falcon 9 launches
- Reset the FlightNumber column
- Replace missing values of PayloadMass with the mean PayloadMass value

1

```
spacex_url="https://api.spacexdata.com/v4/launches/past"
response = requests.get(spacex_url)
```

```
# Use json_normalize method to convert the json result into a dataframe
data = pd.json_normalize(response.json())
```

2

```
# Global variables
BoosterVersion = []
PayloadMass = []
Orbit = []
LaunchSite = []
Outcome = []
Flights = []
GridFins = []
Reused = []
Legs = []
LandingPad = []
Block = []
ReusedCount = []
Serial = []
Longitude = []
Latitude = []
```

```
# Call getBoosterVersion
getBoosterVersion(data)
```

```
# Call getLaunchSite
getLaunchSite(data)
```

```
# Call getPayloadData
getPayloadData(data)
```

```
# Call getCoreData
getCoreData(data)
```

```
launch_dict = {'FlightNumber': list(data['flight_number']),
'Date': list(data['date']),
'BoosterVersion':BoosterVersion,
'PayloadMass':PayloadMass,
'Orbit':Orbit,
'LaunchSite':LaunchSite,
'Outcome':Outcome,
'Flights':Flights,
'GridFins':GridFins,
'Reused':Reused,
'Legs':Legs,
'LandingPad':LandingPad,
'Block':Block,
'ReusedCount':ReusedCount,
'Serial':Serial,
'Longitude': Longitude,
'Latitude': Latitude}
```

3

```
# Create a data from launch_dict
df = pd.DataFrame.from_dict(launch_dict)
```

4

```
data_falcon9 = df[df['BoosterVersion']!='Falcon 1']
```

```
data_falcon9.loc[:, 'FlightNumber'] = list(range(1, data_falcon9.shape[0]+1))
```

```
# Calculate the mean value of PayloadMass column and Replace the np.nan values with its mean value
data_falcon9 = data_falcon9.fillna(value={'PayloadMass': data_falcon9['PayloadMass'].mean()})
```

# DATA COLLECTION – WEB SCRAPING

Web scraping to collect Falcon 9 historical launch records from a Wikipedia page titled List of Falcon 9 and Falcon Heavy launches.

- 1
- Request the HTML page from the static URL
  - Assign the response to an object

- 2
- Create a BeautifulSoup object from the HTML response object
  - Find all tables within the HTML page

- 3
- Collect all column header names from the tables found within the HTML page

- 4
- Use the column names as keys in a dictionary
  - Use custom functions and logic to parse all launch tables (see Appendix) to fill the dictionary values

- 5
- Convert the dictionary to a Pandas DataFrame ready for export

1

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

# use requests.get() method with the provided static_url
response = requests.get(static_url)
# assign the response to a object
data = response.text
```

2

```
soup = BeautifulSoup(data, 'html5lib')
html_tables = soup.find_all('table')
```

3

```
column_names = []

# Apply find_all() function with `th` element on first_launch_table
# Iterate each th element and apply the provided extract_column_from_header() to get a column name
# Append the Non-empty column name ('if name is not None and Len(name) > 0') into a list called column_names

for row in first_launch_table.find_all('th'):
    name = extract_column_from_header(row)
    if(name != None and len(name) > 0):
        column_names.append(name)
```

4

```
launch_dict= dict.fromkeys(column_names)

# Remove an irrelevant column
del launch_dict['Date and time ( )']

# Let's initial the Launch_dict with each value to be an empty list
launch_dict['Flight No.'] = []
launch_dict['Launch site'] = []
launch_dict['Payload'] = []
launch_dict['Payload mass'] = []
launch_dict['Orbit'] = []
launch_dict['Customer'] = []
launch_dict['Launch outcome'] = []
# Added some new columns
launch_dict['Version Booster']=[]
launch_dict['Booster landing']=[]
launch_dict['Date']= []
launch_dict['Time']= []
```

# 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
  - 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.
  - Create a Class column that contains the values from the list landing\_class
  - Export the DataFrame as a .csv file.

1

```
bad_outcomes=set(landing_outcomes.keys()|[1,3,5,6,7])  
bad_outcomes  
['False ASDS', 'False Ocean', 'False RTLS', 'None ASDS', 'None None']
```

3

```
# Landing_class = 0 if bad_outcome  
# Landing_class = 1 otherwise  
  
landing_class = []  
  
for outcome in df['Outcome']:  
    if outcome in bad_outcomes:  
        landing_class.append(0)  
    else:  
        landing_class.append(1)
```

3

```
df['Class']=landing_class
```

4

```
df.to_csv("dataset_part\2.csv", index=False)
```

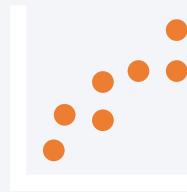
# EXPLORATORY DATA ANALYSIS (EDA) - VISUALIZATION

---

## SCATTER CHARTS

Scatter charts were produced to visualize the relationships between:

- Flight Number and Launch Site
- Payload and Launch Site
- Orbit Type and Flight Number
- Payload and Orbit Type



Scatter charts are useful to observe relationships, or correlations, between two numeric variables.

## BAR CHART

A bar chart was produced to visualize the relationship between:

- Success Rate and Orbit Type

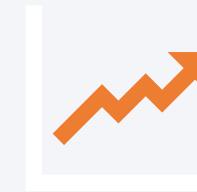


Bar charts are used to compare a numerical value to a categorical variable. Horizontal or vertical bar charts can be used, depending on the size of the data.

## LINE CHARTS

Line charts were produced to visualize the relationships between:

- Success Rate and Year (i.e. the launch success yearly trend)



Line charts contain numerical values on both axes, and are generally used to show the change of a variable over time.

# EXPLORATORY DATA ANALYSIS (EDA) – SQL

---

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

The SQL queries 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
8. 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

# INTERACTIVE DASHBOARD – PLOTLY DASH

---

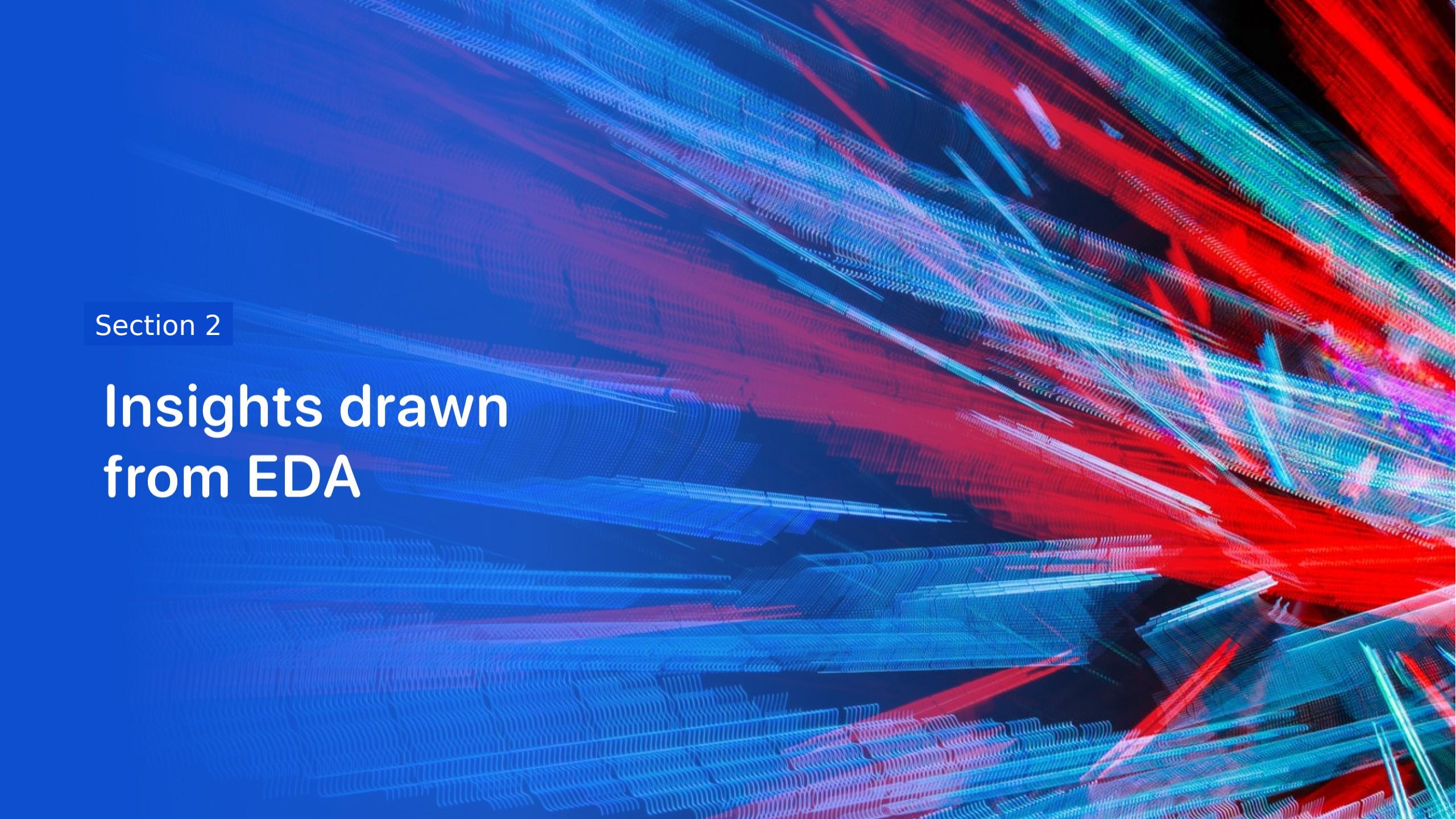
The following plots were added to a Plotly Dash dashboard to have an interactive visualisation of the data:

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

# Results

---

- Exploratory data analysis results
- Interactive analytics demo in screenshots
- Predictive analysis results

The background of the slide features a complex, abstract pattern of wavy, glowing lines in shades of blue, red, green, and purple. These lines are arranged in several parallel bands that curve and twist across the frame, creating a sense of depth and motion.

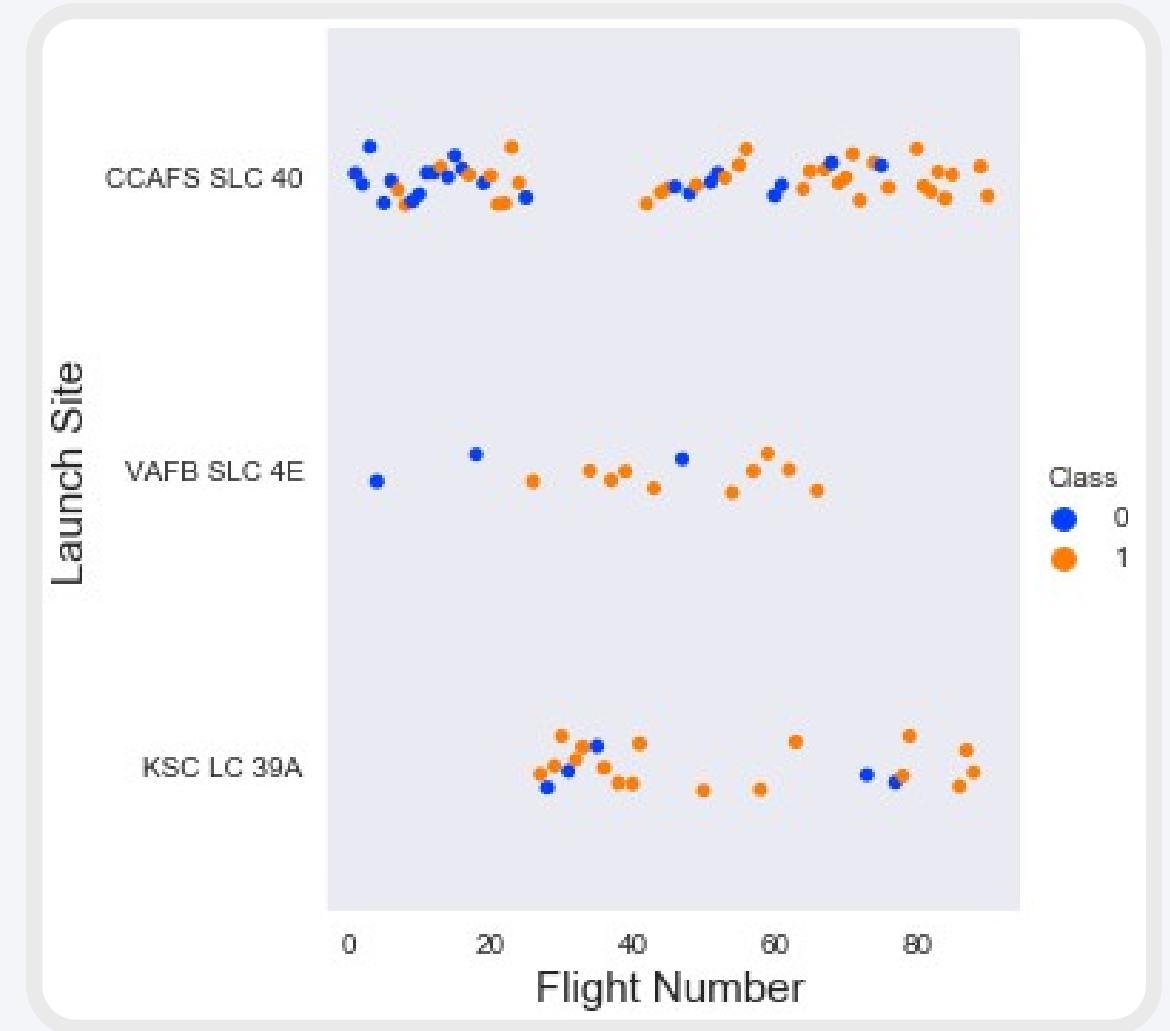
Section 2

## Insights drawn from EDA

# LAUNCH SITE VS. FLIGHT NUMBER

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

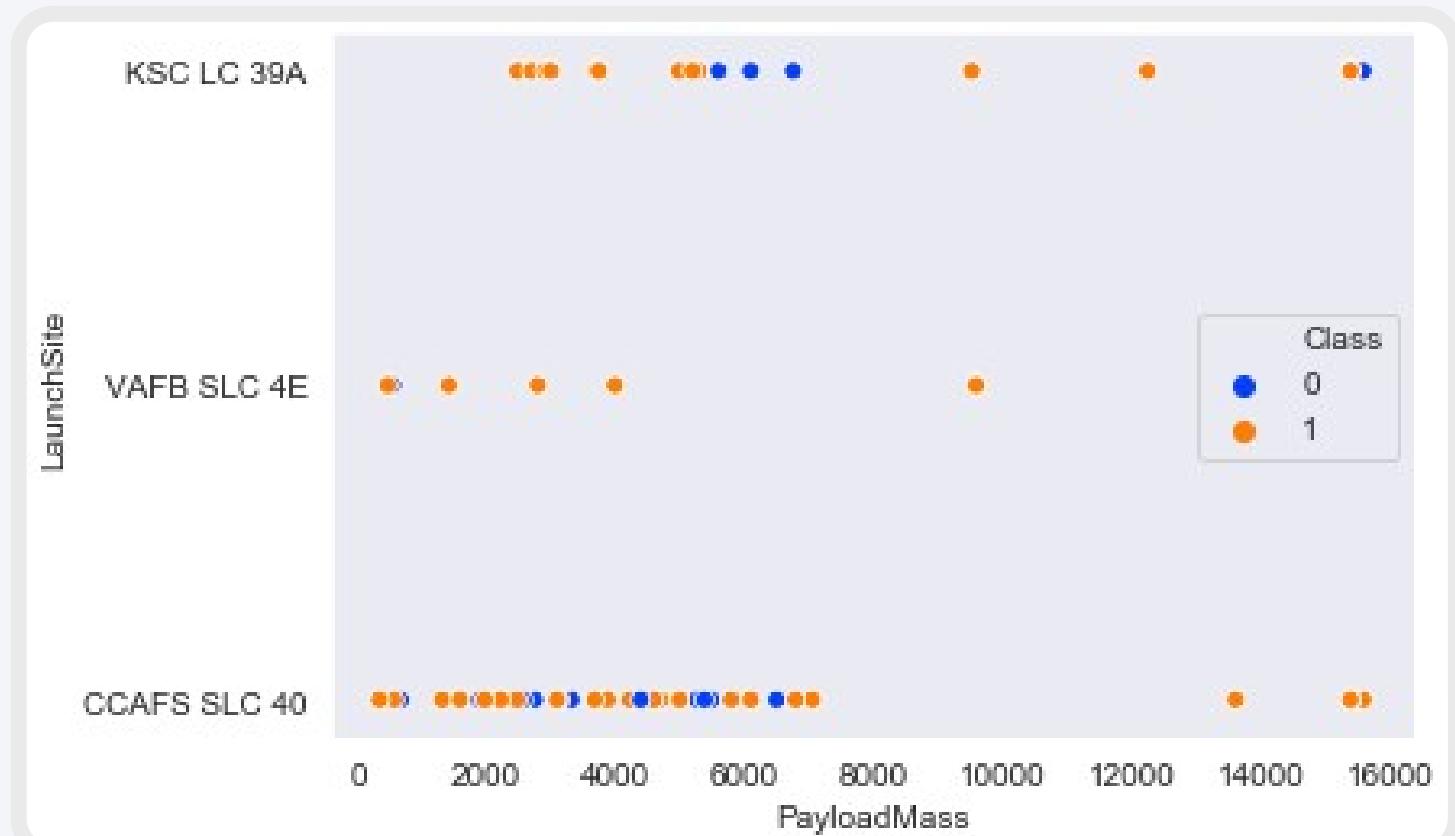
- As the number of flights increases, the rate of success at a launch site increases.
- Most of the early flights (flight numbers < 30) were launched from CCAFS SLC 40, and were generally unsuccessful.
- The flights from VAFB SLC 4E also show this trend, that earlier flights were less successful.
- No 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).



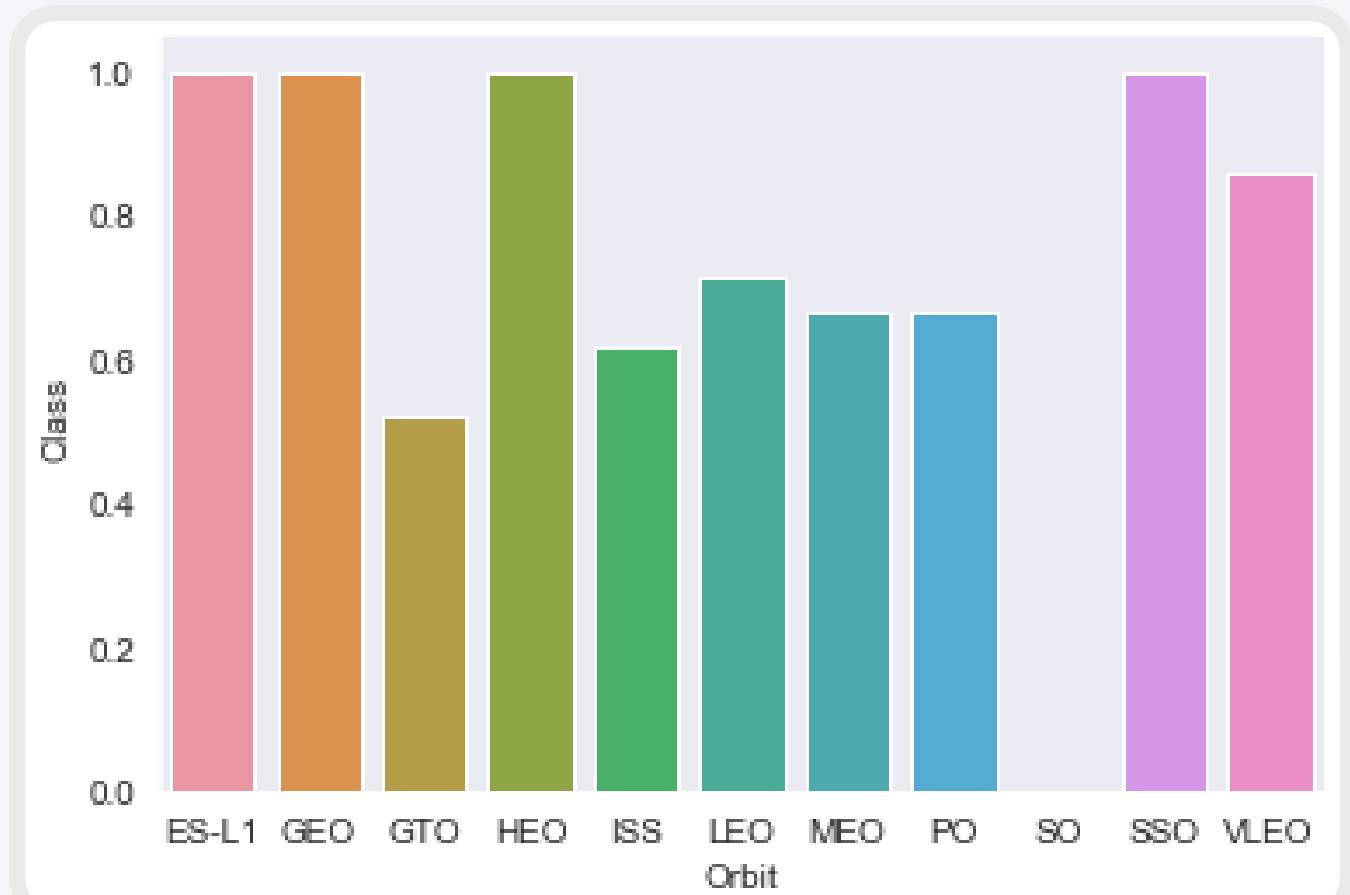
# SUCCESS RATE VS. ORBIT TYPE

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:

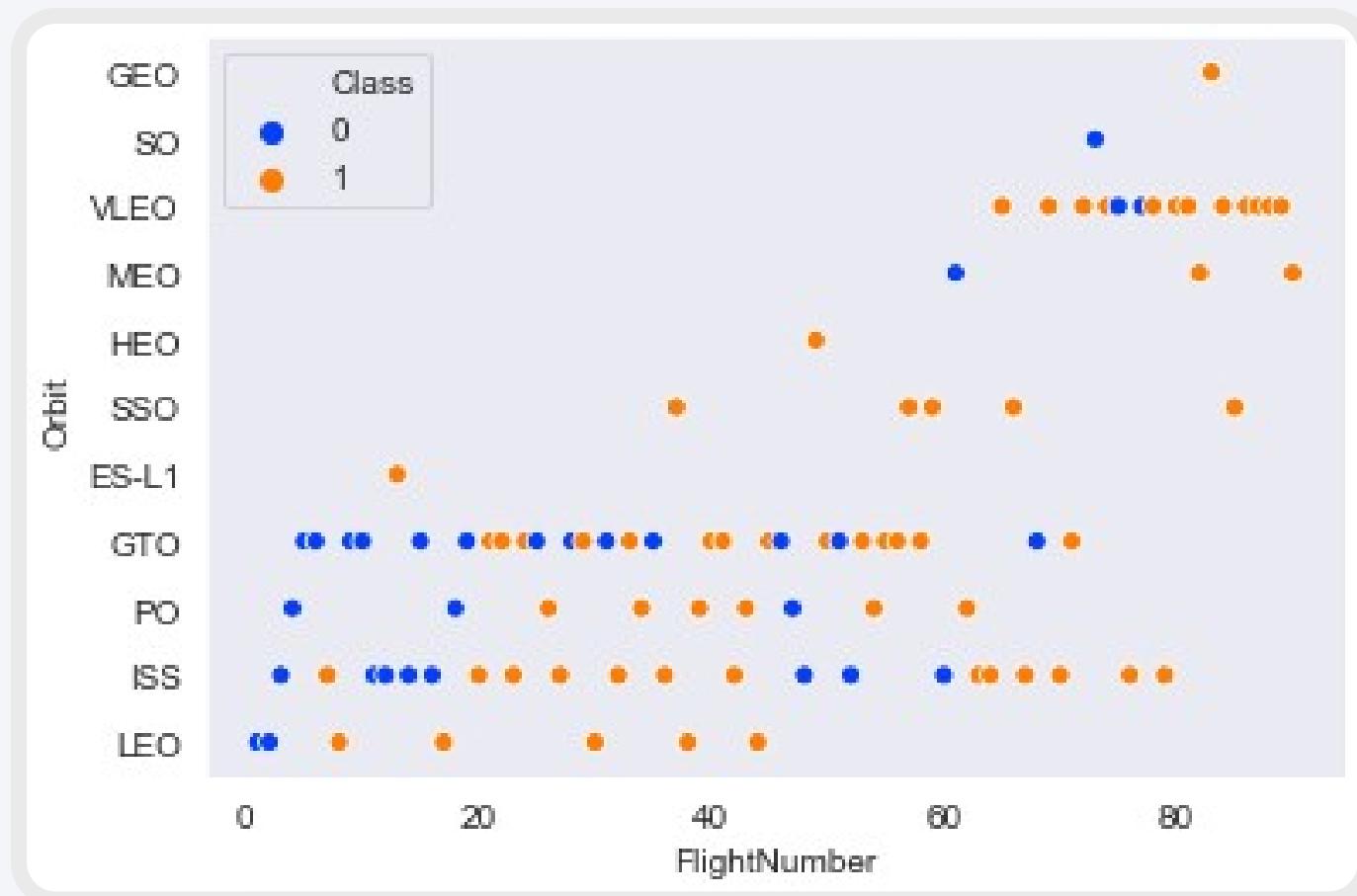
- SO (Heliocentric Orbit)



# 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:

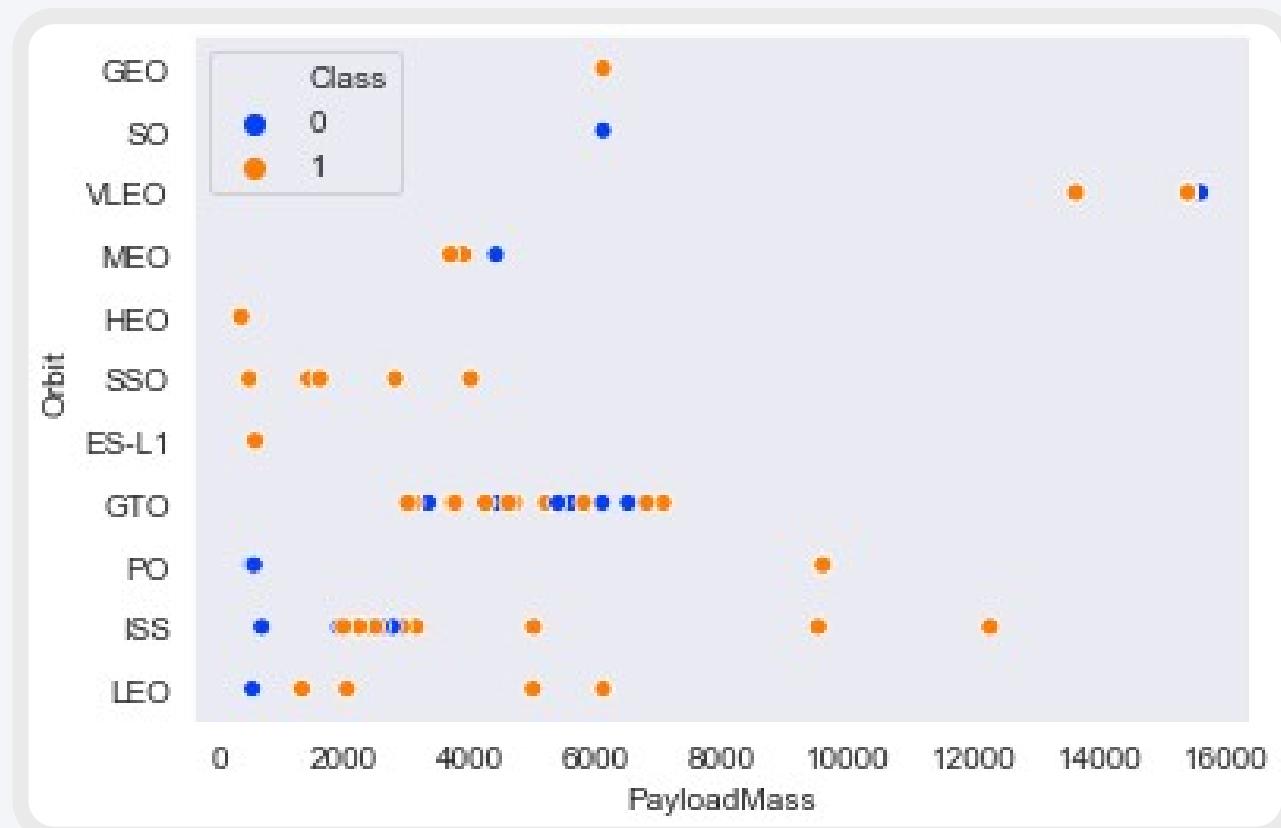
- The 100% success rate of GEO, HEO, and ES-L1 orbits can be explained by only having 1 flight into the respective orbits.
- The 100% success rate in SSO is more impressive, with 5 successful flights.
- There is little relationship between Flight Number and Success Rate for GTO.
- Generally, as Flight 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:

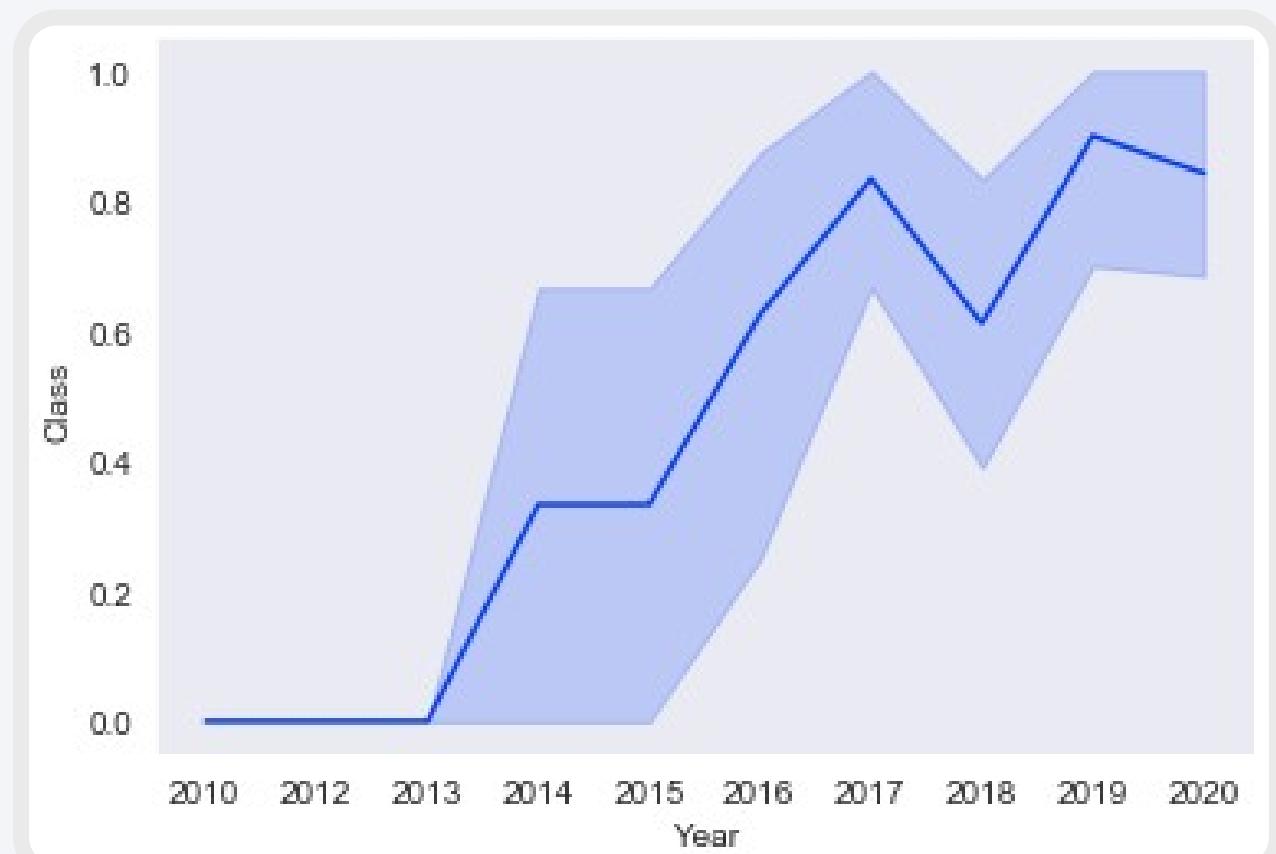
- The following orbit types have more success with heavy payloads:
  - PO (although the number of data points is small)
  - ISS
  - LEO
- For GTO, the relationship between payload mass and success rate is unclear.
- VLEO (Very Low Earth Orbit) launches are associated with heavier payloads, which makes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



# All Launch Site Names

---

```
%sql select unique(launch_site) from SPACEXTBL
```

The word UNIQUE returns only unique values from the LAUNCH\_SITE column of the SPACEXTBL table.

# Launch Site Names Begin with 'CCA'

---

```
%sql select lunch_Site from SPACEXTBL where lunch_site like  
`CCA%` limit 5;
```

LIMIT 5 fetches only 5 records, and the LIKE keyword is used with the wild card 'CCA%' to retrieve string values beginning with 'CCA'.

# Total Payload Mass

---

```
%sq select sum(PAYLOAD_MASS_KG_) as total_payload from  
SPACEXTBL;
```

455966 is out put

The SUM keyword is used to calculate the total of the LAUNCH column, and the SUM keyword (and the associated condition) filters the results to only boosters from NASA (CRS).

# Average Payload Mass by F9 v1.1

---

```
%sq select avg(PAYLOAD_MASS_KG_) as avg_payload from  
SPACEXTBL where BOOSTER_VERSION='F9 v1.1';
```

2928 is output

The AVG keyword is used to calculate the average of the PAYLOAD\_MASS\_KG\_ column, and the WHERE keyword (and the associated condition) filters the results to only the F9 v1.1 booster version.

# First Successful Ground Landing Date

---

```
%sql select min(Date) from SPACEXTABLE where  
Landing_Outcome like '%success%'
```

The MIN keyword is used to calculate the minimum of the DATE column, i.e. the first date, and the WHERE keyword (and the associated condition) filters the results to only the successful ground pad landings.

## Successful Drone Ship Landing with Payload between 4000 and 6000

---

List the names of boosters which have successfully landed on drone ship and had payload mass greater than 4000 but less than 6000

```
%sql select Booster_Version from SPACEXTABLE where  
PAYLOAD_MASS_KG_ between 4000 and 6000
```

The WHERE keyword is used to filter the results to include only those that satisfy both conditions in the brackets (as the AND keyword is also used). The BETWEEN keyword allows for  $4000 < x < 6000$  values to be selected.

# Total Number of Successful and Failure Mission Outcomes

---

Calculate the total number of successful and failure mission outcomes

```
%sql select Mission_Outcome,count(Mission_Outcome) from  
SPACEXTABLE group by Mission_Outcome
```

```
Mission_Outcome      count(Mission_Outcome)
```

```
Failure (in flight) 1
```

```
Success   98
```

```
Success   1
```

```
Success (payload status unclear) 1
```

The COUNT keyword is used to calculate the total number of mission outcomes, and the GROUPBY keyword is also used to group these results by the type of mission outcome.

# Boosters Carried Maximum Payload

---

List the names of the booster which have carried the maximum payload mass

```
%sql select Booster_Version,PAYLOAD_MASS_KG_ from SPACEXTABLE where PAYLOAD_MASS_KG_ >= (select max(PAYLOAD_MASS_KG_) from SPACEXTABLE )
```

Booster_Version	PAYLOAD_MASS_KG_
F9 B5 B1048.4	15600
F9 B5 B1049.4	15600
F9 B5 B1051.3	15600
F9 B5 B1056.4	15600
F9 B5 B1048.5	15600

A subquery is used here. The SELECT statement within the brackets finds the maximum payload, and this value is used in the WHERE condition. The DISTINCT keyword is then used to retrieve only distinct /unique booster versions.

# 2015 Launch Records

---

```
%sql select substr(Date, 6,2) as month, Landing_Outcome ,  
Booster_Version from SPACEXTABLE where substr(Date,  
0,5)='2015' and Landing_Outcome like '%Failure%'
```

The WHERE keyword is used to filter the results for only failed landing outcomes, AND only for the year of 2015.

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

---

```
%sql select Date,Landing_Outcome, count(Landing_Outcome) from  
SPACEXTABLE group by Landing_Outcome having Date between  
'2010-06-04' and '2017-03-20'
```

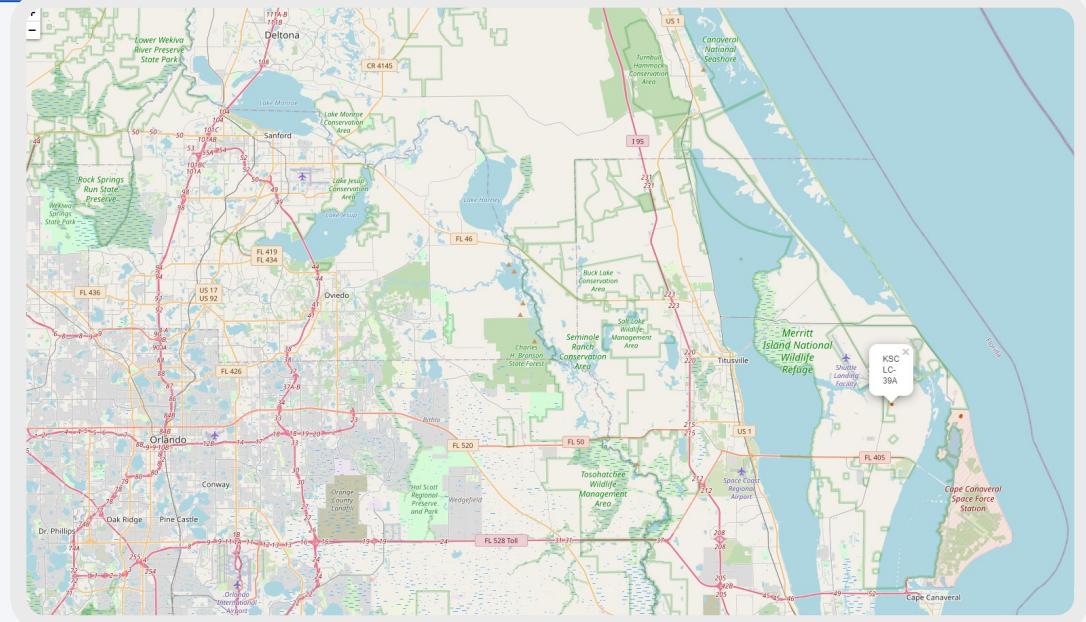
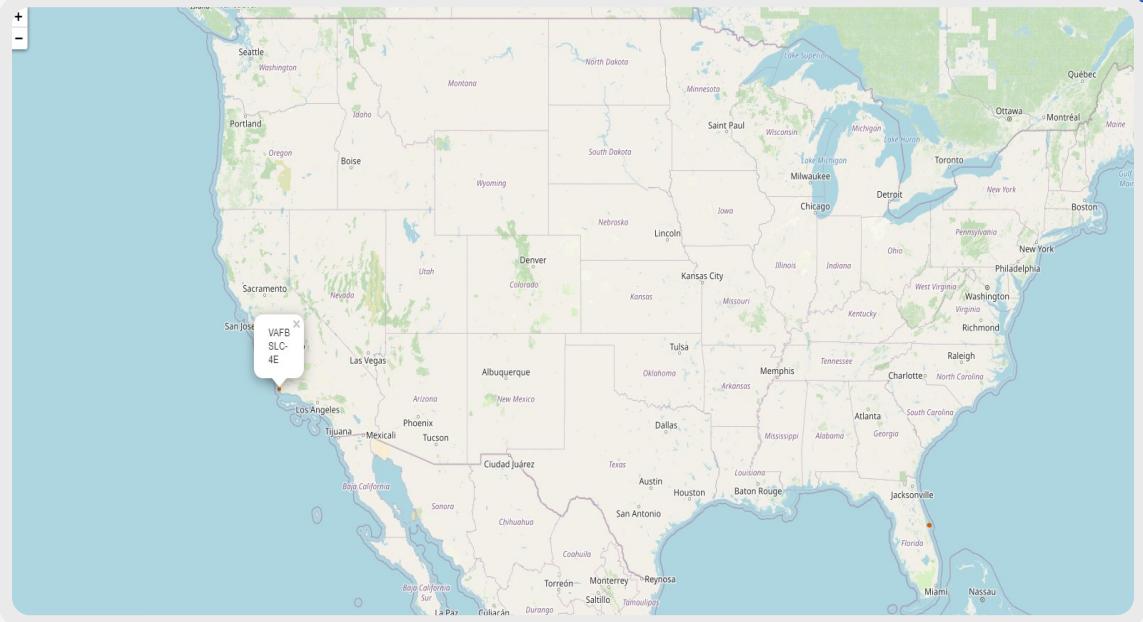
The WHERE keyword is used with the BETWEEN keyword to filter the results to dates only within those specified. The results are then grouped and ordered, using the keywords GROUP BY and ORDER BY, respectively, where DESC is used to specify the descending order.

The background of the slide is a photograph taken from space at night. It shows the curvature of the Earth's horizon against a dark blue sky. City lights are visible as numerous small white and yellow dots, primarily concentrated in the lower right quadrant where the United States appears. In the upper right, there are bright green and yellow bands of light, likely the Aurora Borealis or Australis. The overall atmosphere is dark and mysterious.

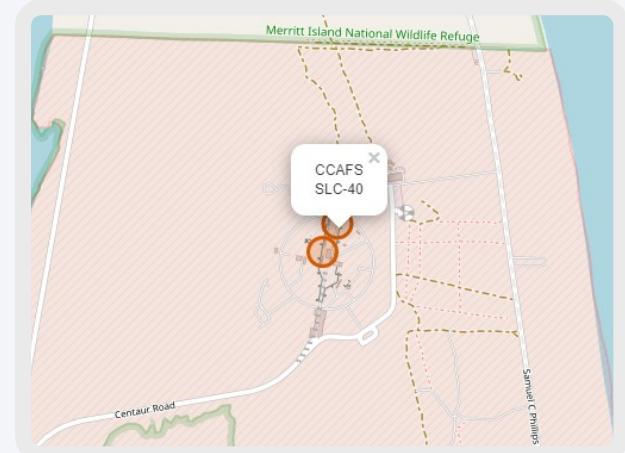
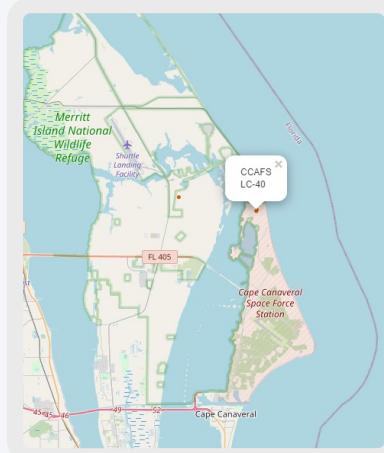
Section 3

# Launch Sites Proximities Analysis

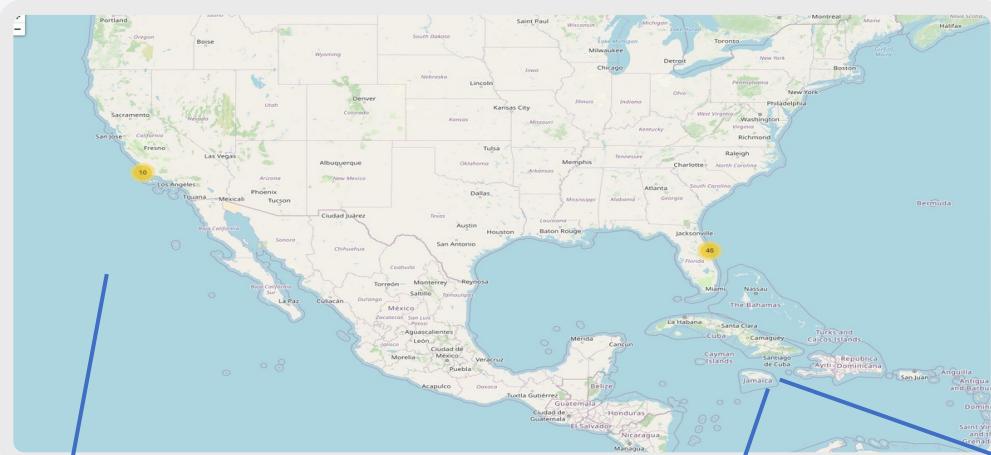
# ALL LAUNCH SITES ON A MAP



All SpaceX launch sites are on coasts of the United States of America, specifically Florida and California.



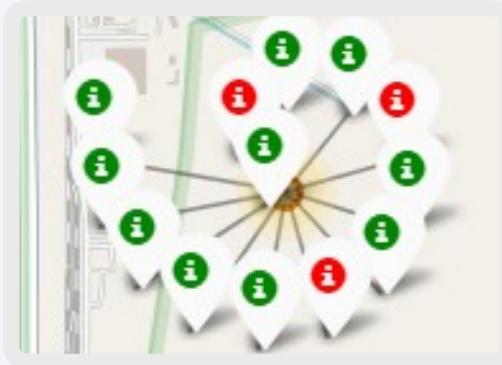
# SUCCESS/FAILED LAUNCHES FOR EACH SITE



VAFB SLC-4E

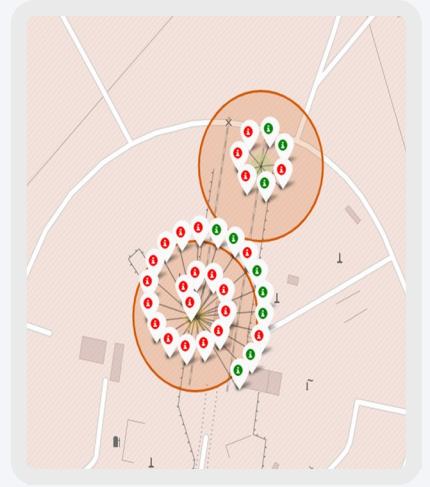


KSC LC-39A



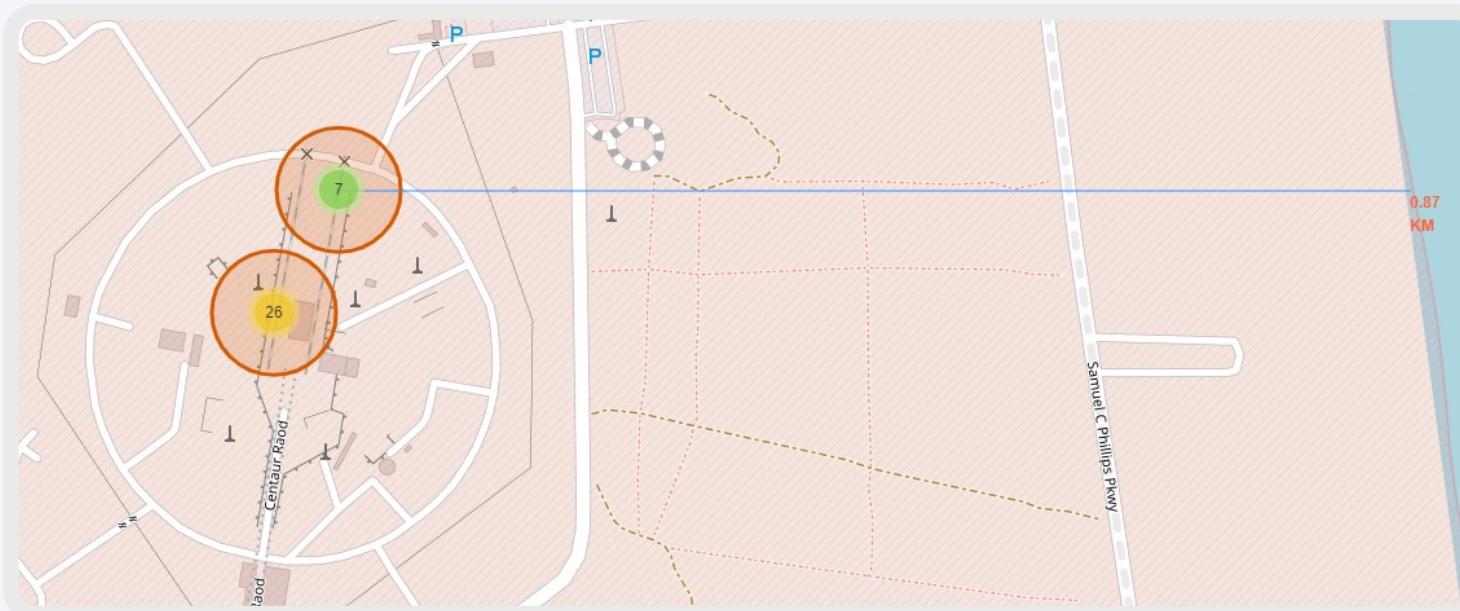
Launches have been grouped into clusters, and annotated with green icons for successful launches, and red icons for failed launches.

CCAFS SLC-40 and CCAFS LC-40



# PROXIMITY OF LAUNCH SITES TO OTHER POINTS OF INTEREST

Using the CCAFS SLC-40 launch site as an example site, we can understand more about the placement of launch sites.



Are launch sites in close proximity to railways?

- YES. The coastline is only 0.87 km due East.

Are launch sites in close proximity to highways?

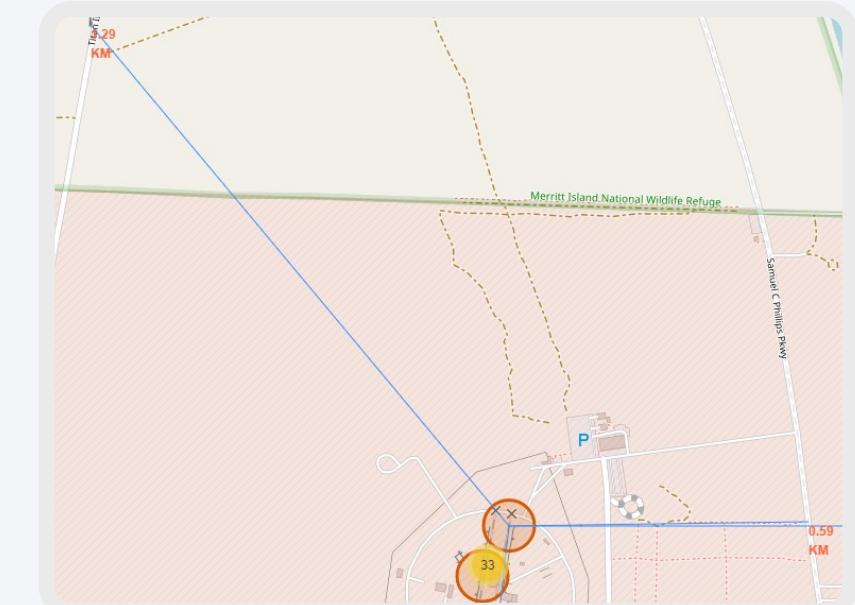
- YES. The nearest highway is only 0.59km away.

Are launch sites in close proximity to railways?

- YES. The nearest railway is only 1.29 km away.

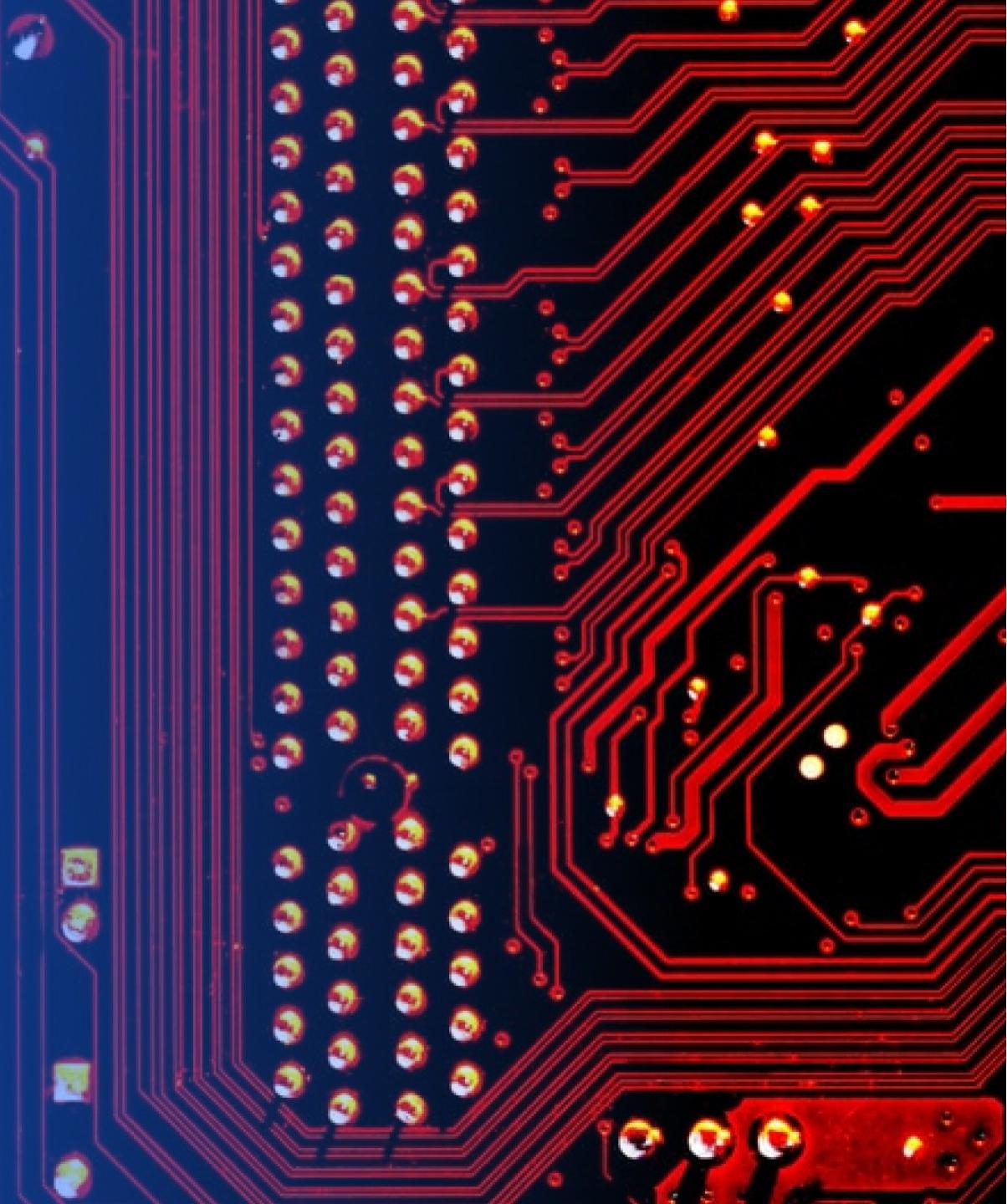
Do launch sites keep certain distance away from cities?

- YES. The nearest city is 51.74 km away.

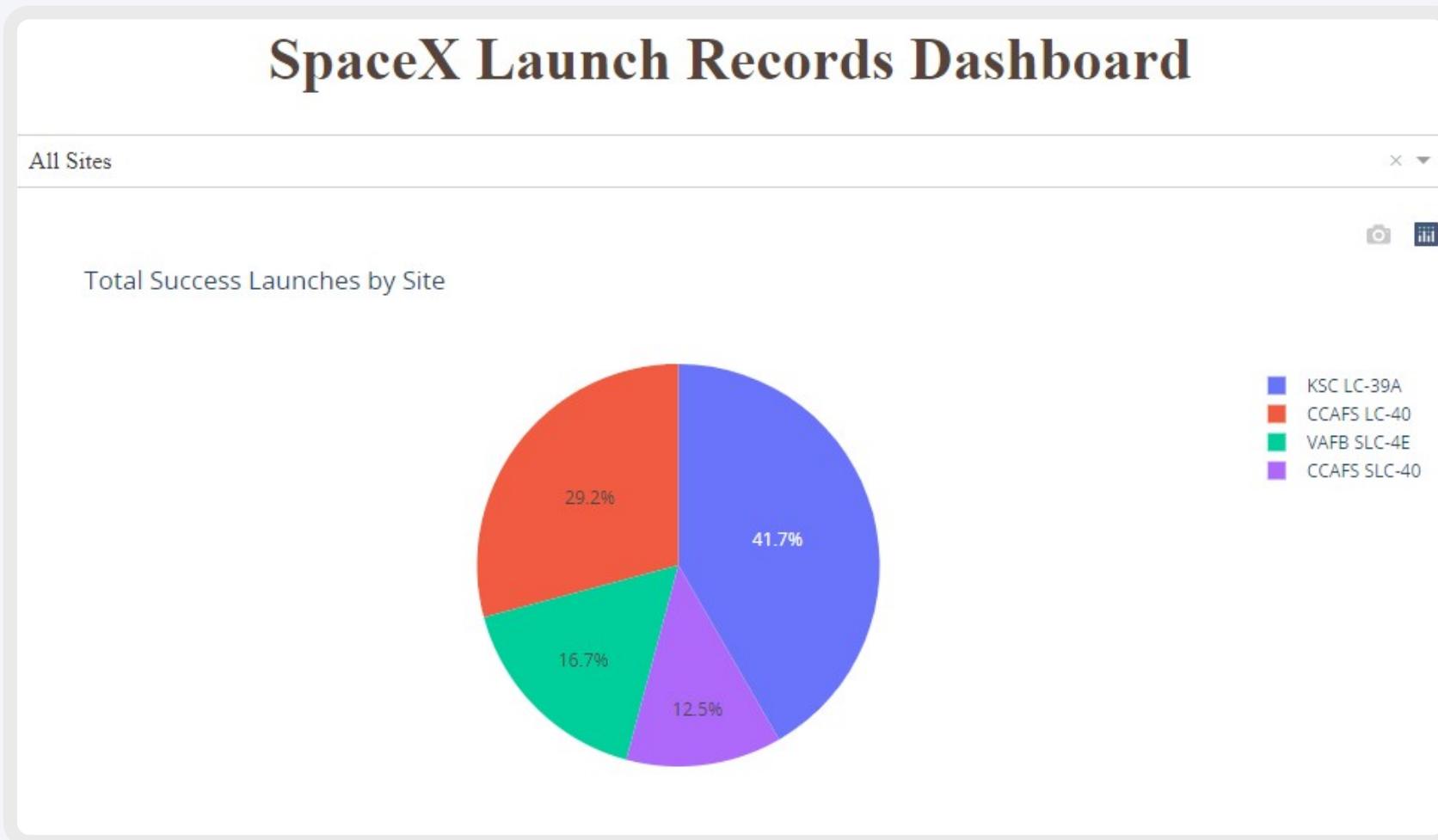


Section 4

# Build a Dashboard with Plotly Dash

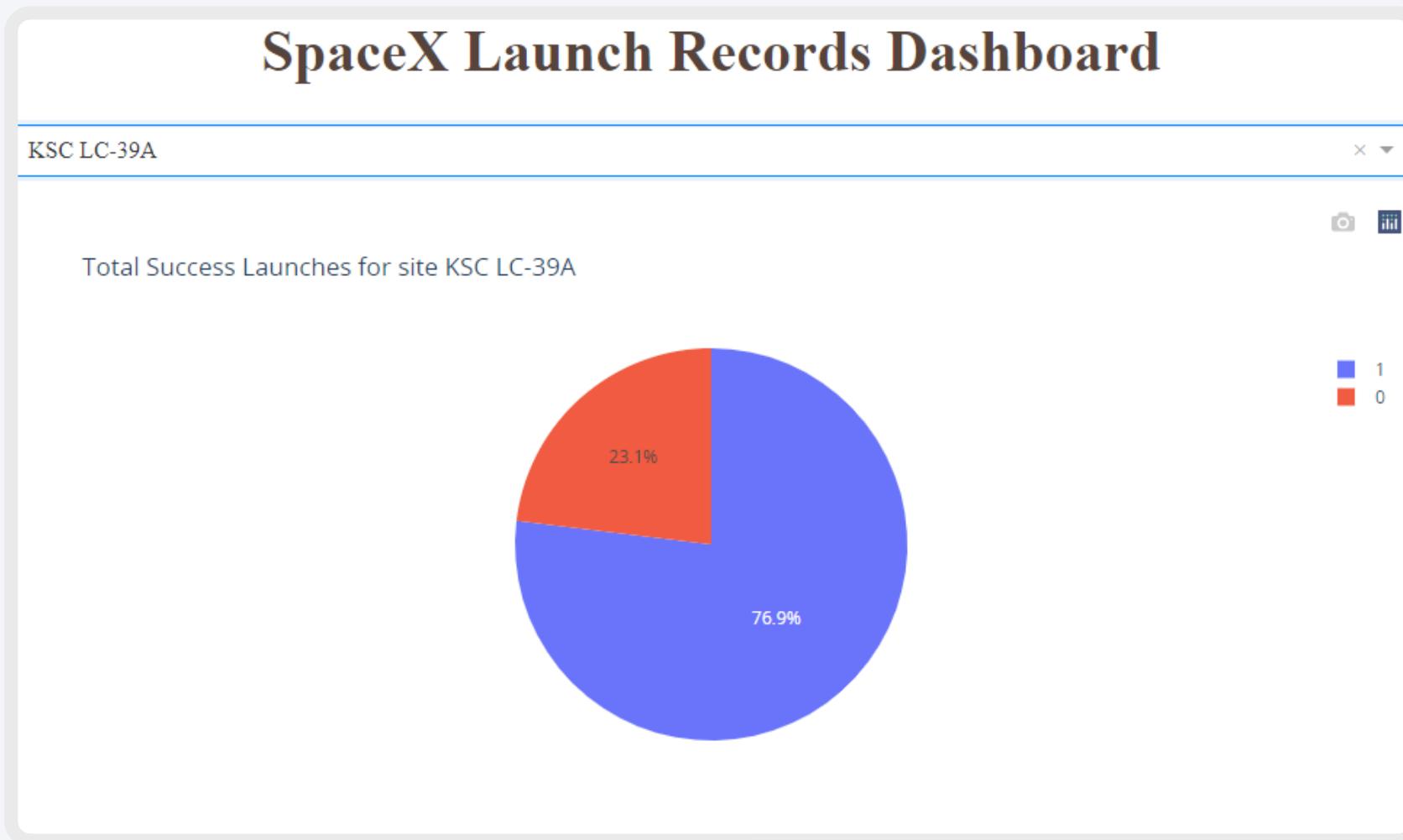


# LAUNCH SUCCESS COUNT FOR ALL SITES



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

# PIE CHART FOR THE LAUNCH SITE WITH HIGHEST LAUNCH SUCCESS RATIO



The launch site **KSC LC-39 A** also had the highest rate of successful launches, with a 76.9% success rate.

# LAUNCH OUTCOME VS. PAYLOAD SCATTER PLOT FOR ALL SITES



- Plotting the launch outcome vs. payload for all sites shows a gap around 4000 kg, so it makes sense to split the data into 2 ranges:
  - 0 – 4000 kg (low payloads)
  - 4000 – 10000 kg (massive payloads)
- From these 2 plots, it can be shown that the success for massive payloads is lower than that for low payloads.
- It is also worth noting that some booster types (v1.0 and B5) have not been launched with massive payloads.

Section 5

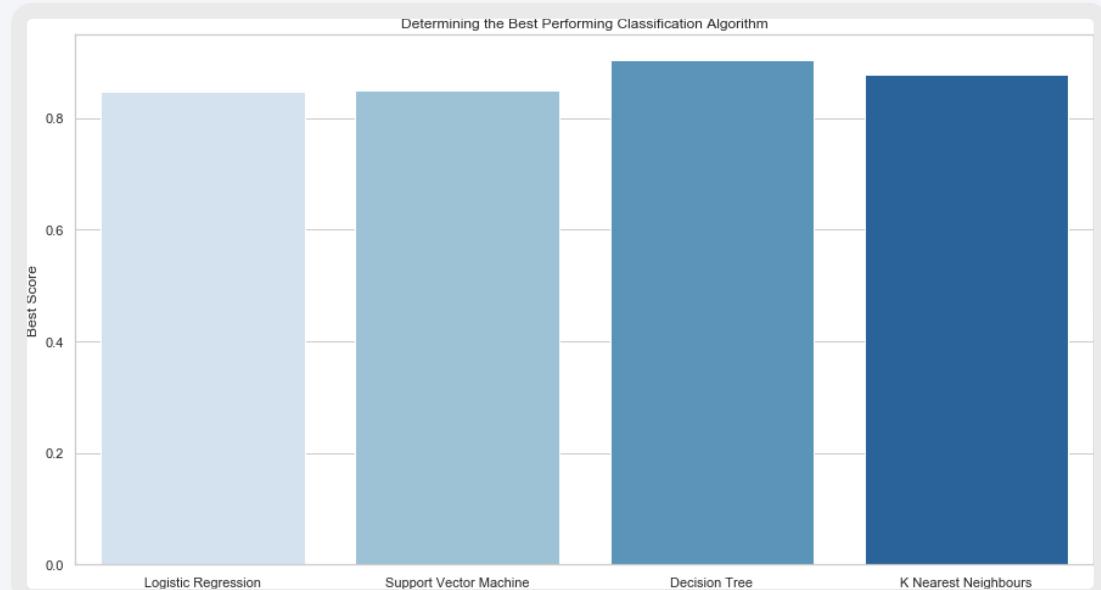
# Predictive Analysis (Classification)

# CLASSIFICATION ACCURACY

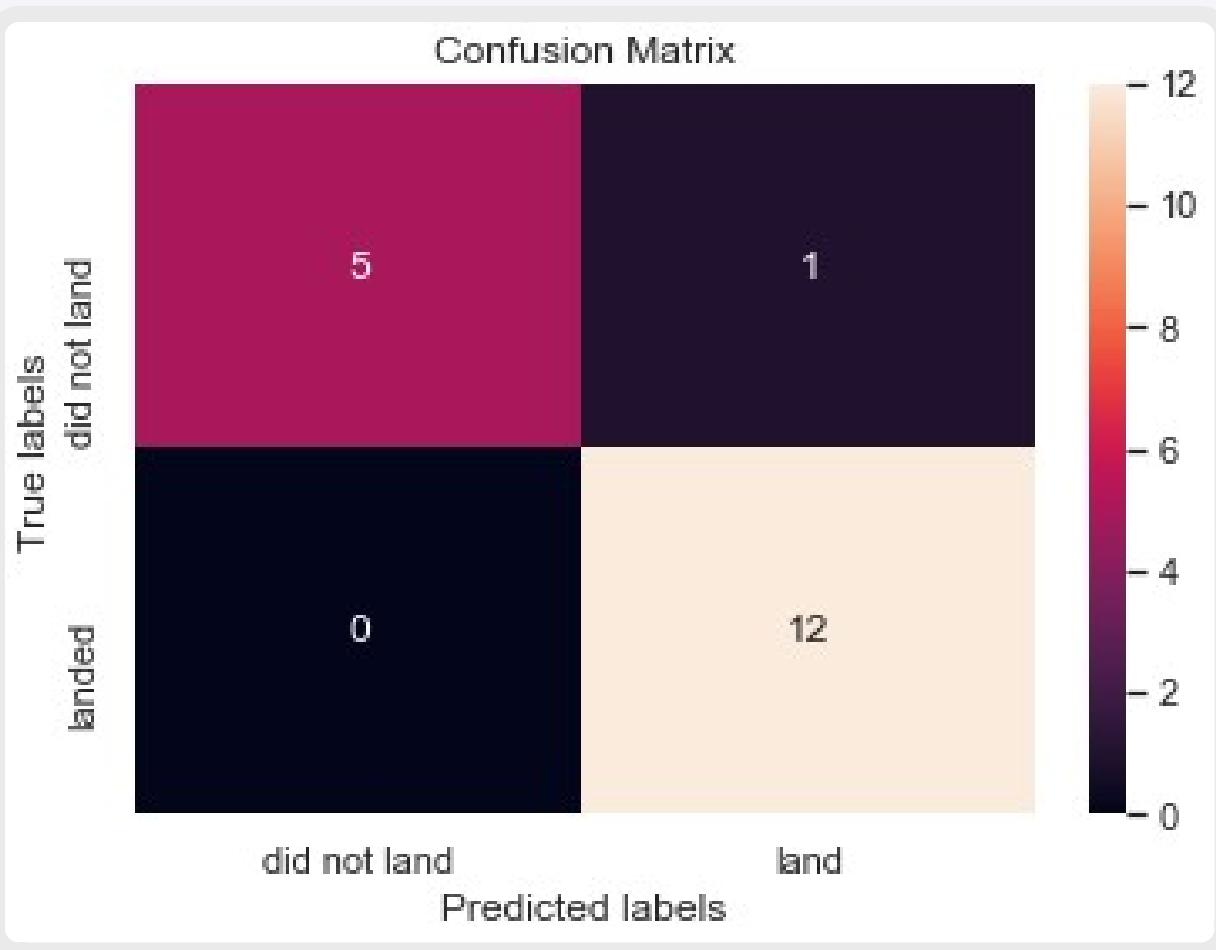
Plotting the Accuracy Score and Best Score for each classification algorithm produces the following result:

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

Algorithm	Accuracy Score	Best Score
Logistic Regression	0.833333	0.846429
Support Vector Machine	0.833333	0.848214
Decision Tree	0.944444	0.903571
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

---

- 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%.

# Appendix

---

Custom functions to retrieve the required information

Custom logic to clean the data

Thank you!

