



IBM Developer  
SKILLS NETWORK

# Winning Space Race with Data Science

JABRANE MERIZAK  
JULY 2023



# 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 resul

# Introduction

---

- **Project background and context**

SpaceX advertises Falcon 9 rocket launches on its website with a cost of 62 million dollars; other providers cost upwards of 165 million dollars each, much of the savings is because SpaceX can reuse the first stage. Therefore, if we can determine if the first stage will land, we can determine the cost of a launch. SpaceX's Falcon 9 launch like regular rockets.

- **Problems you want to find answers**

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

---

## Executive Summary

- Data collection methodology:
  - Data was collected using SpaceX API and Webscraping from Wikipedia.
- Perform data wrangling
  - Data was processed by converting categorical variables and one-hot encoding.
- 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 collection process followed these steps:**
  - 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

---

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

```
[7]: spacex_url="https://api.spacexdata.com/v4/launches/past"
```

```
[9]: response = requests.get(spacex_url)
```

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

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



# Data Collection - Scrapping

---

```
: # use requests.get() method with the provided static_url  
# assign the response to a object  
response = requests.get(static_url).text  
response[:100]
```

```
# Use BeautifulSoup() to create a BeautifulSoup object from  
soup = BeautifulSoup(response)
```

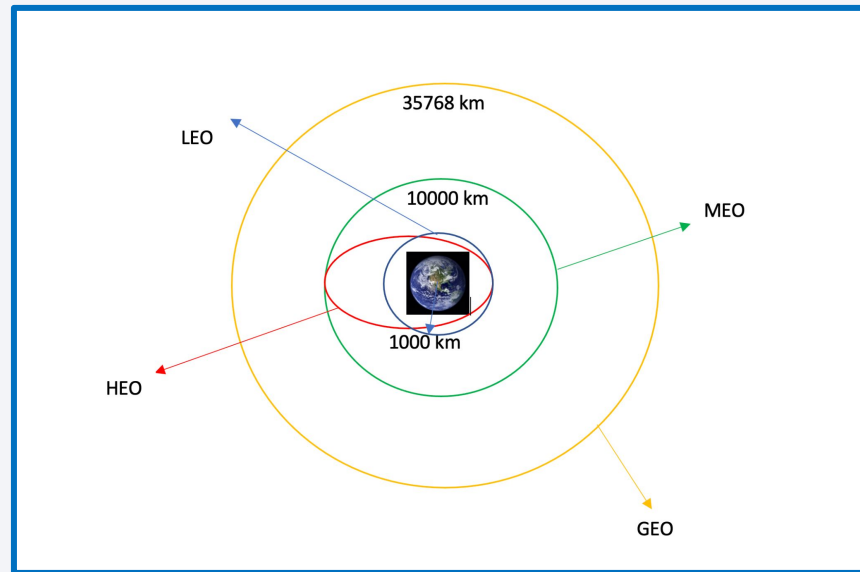
```
html_pd_tables = pd.read_html(static_url)  
html_pd_tables[2].head()
```

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

# Data Wrangling

---

- We performed exploratory data analysis and determined the training labels.
- We calculated the number of launches at each site, and the number and occurrence of each orbits
- We created landing outcome label from
- outcome column and exported the results to csv.



# EDA with Data 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

## LINE CHARTS

Line charts were produced to visualize the relationships between:

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

## BAR CHART

A bar chart was produced to visualize the

relationship between:

- Success Rate and Orbit Type

# EDA with 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

# 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.

# 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.

# 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.



# Results

---

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

The background of the slide is an abstract composition. It features a dark blue field on the left side, which transitions into a complex pattern of diagonal streaks in vibrant red and cyan. These streaks are layered over a fine, light blue grid that covers the entire right half of the image. The overall effect is one of dynamic energy and technological sophistication.

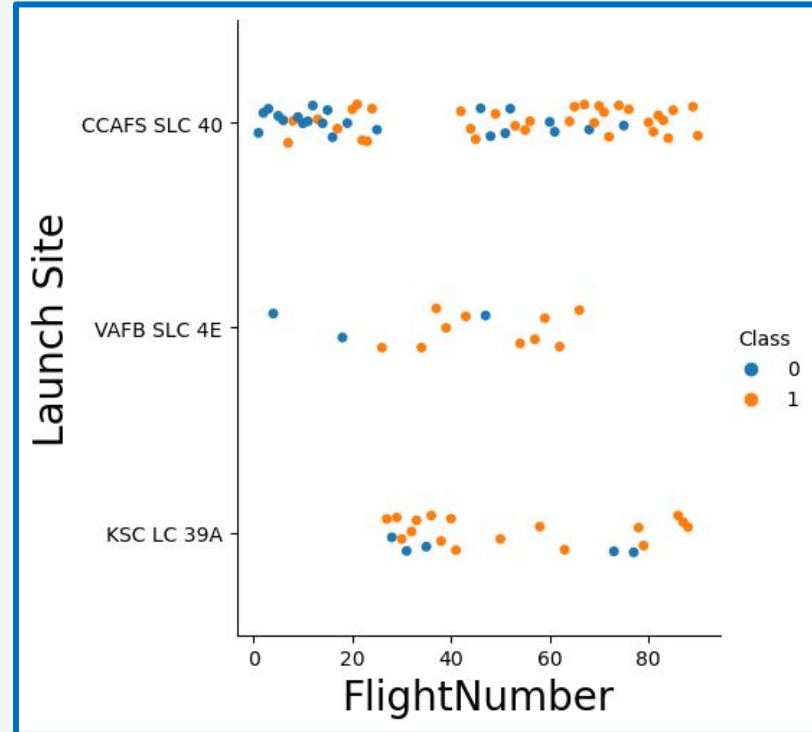
Section

2

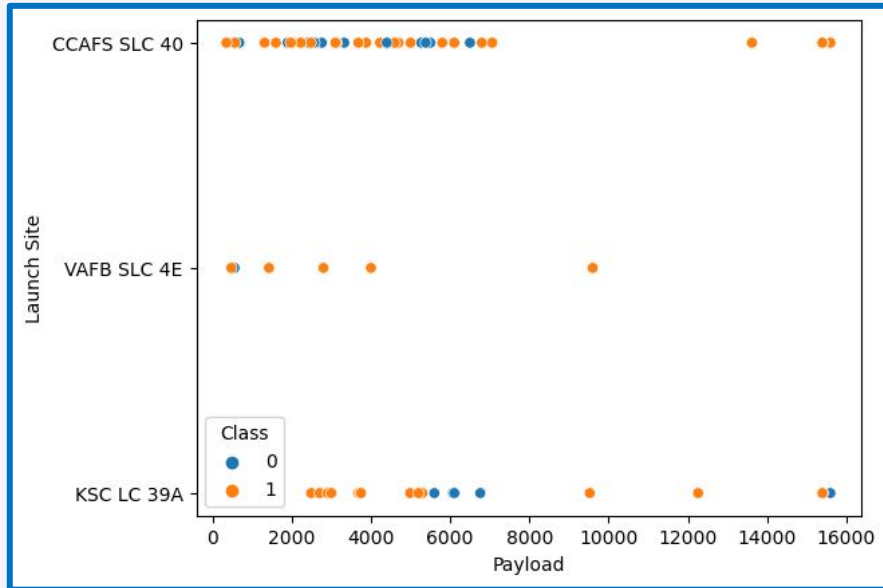
# Insights drawn from EDA

# 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



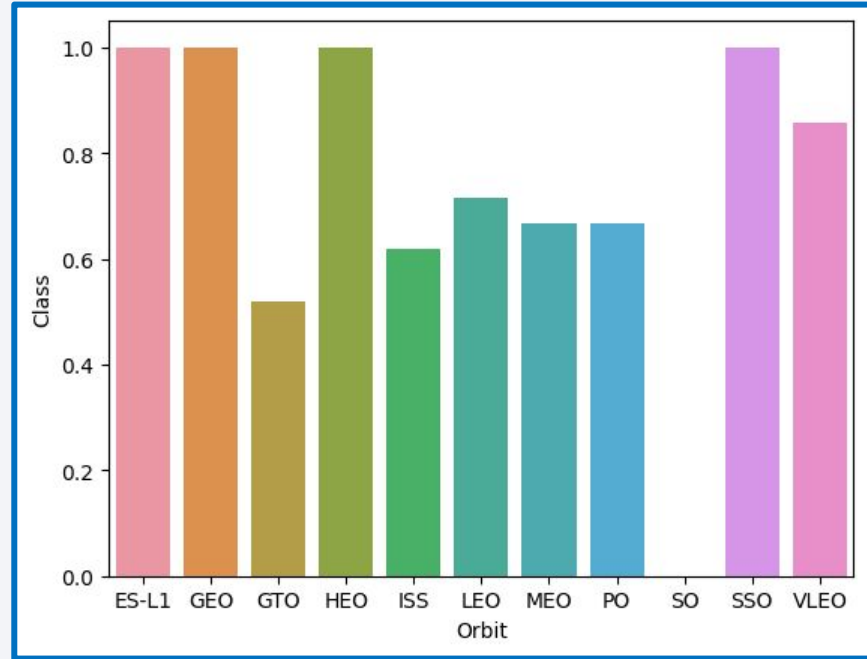
The scatter plot of Launch Site vs. Payload Mass shows that:

- 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

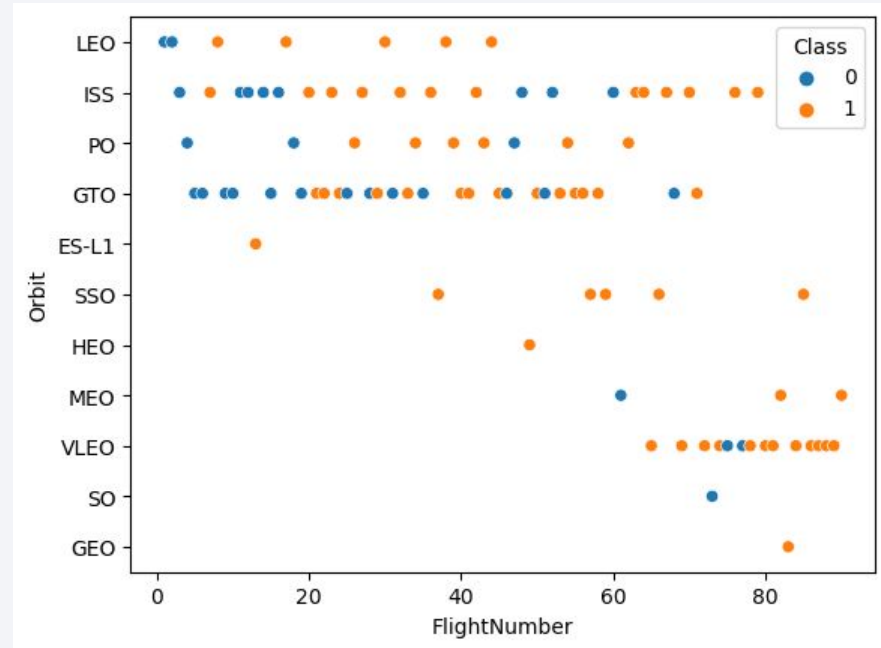
---

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

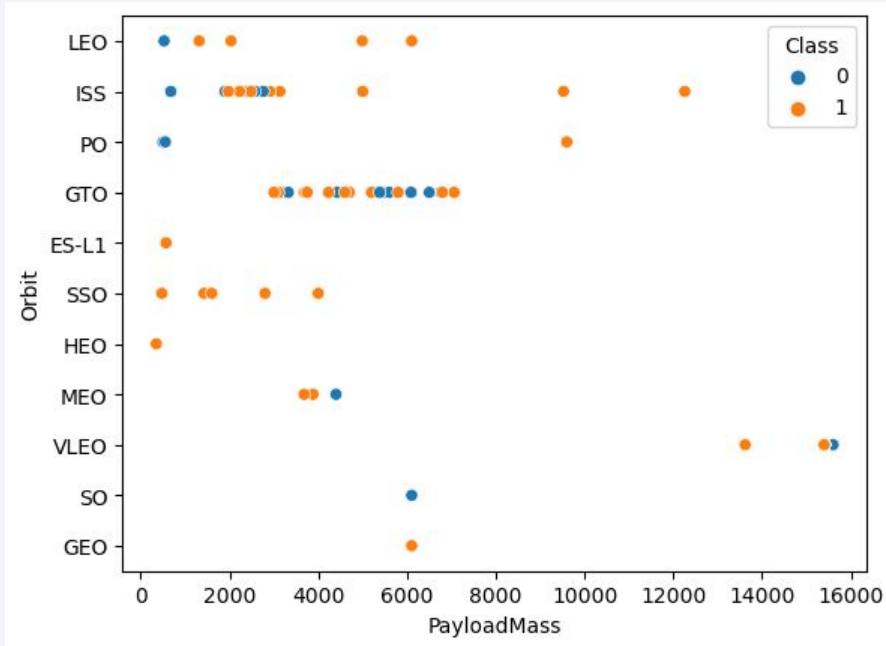


# Flight Number vs. Orbit Type

- 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

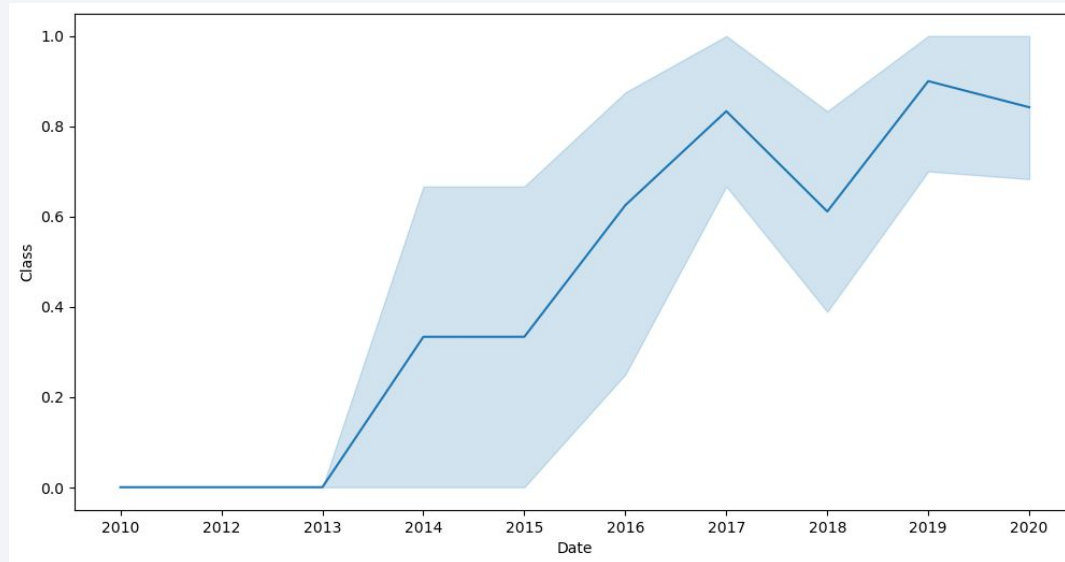


- We can observe that with heavy payloads, the successful landing are more for PO, LEO and ISS orbits



# 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 key word DISTINCT to show only unique launch sites from the SpaceX data.

```
%sql select distinct(Launch_Site) from spacetable
* sqlite:///my_data1.db
Done.
```

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

# Launch Site Names Begin with 'CCA'

```
%sql select * from spacetable where launch_site like 'CCA%' limit 5
* sqlite:///my_data1.db
Done.
```

Date	Time (UTC)	Booster_Version	Launch_Site	Payload	PAYLOAD_MASS_KG_	Orbit	Customer	Mission_Outcome	Landing_Outcome
06/04/2010	18:45:00	F9 v1.0 B0003	CCAFS LC-40	Dragon Spacecraft Qualification Unit	0.0	LEO	SpaceX	Success	Failure (parachute)
				Dragon demo flight					

- We used the query above to display 5 records where launch sites begin with 'CCA'

# Total Payload Mass

---

- We calculated the total payload carried by boosters from NASA as 45596 using the query below

```
%sql select SUM(PAYLOAD_MASS_KG_) from spacetable where Customer = 'NASA (CRS)'  
* sqlite:///my_data1.db  
Done.  
  
SUM(PAYLOAD_MASS_KG_)  
-----  
45596.0
```

# Average Payload Mass by F9 v1.1

---

```
%sql select AVG(PAYLOAD_MASS_KG_) from spacetable where Booster_Version like 'F9 v1.1%'
* sqlite:///my_data1.db
Done.
AVG(PAYLOAD_MASS_KG_)
2534.6666666666665
```

- We calculated the average payload mass carried by booster version F9 v1.1 as 2928.4

# First Successful Ground Landing Date

---

- We observed that the dates of the first successful landing outcome on ground pad was 22nd December 2015

```
%%sql select Date from spacetable
      where Landing_Outcome = 'Success (ground pad)'
      order by date desc limit 1
```

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

Date
------

22/12/2015
------------

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

---

```
%%sql select Booster_Version from spacetable
      where Landing_Outcome = 'Success (drone ship)'
      and PAYLOAD_MASS_KG_ between 4000 and 6000
```

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

Booster_Version
-----------------

F9 FT B1022
-------------

F9 FT B1026
-------------

F9 FT B1021.2
---------------

F9 FT B1031.2
---------------

- 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

---

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

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

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

Done.

Mission_Outcome	count(Mission_Outcome)
Failure (in flight)	1
Success	99
Success (payload status unclear)	1

# Boosters Carried Maximum Payload

---

```
%%sql select Booster_Version from spacetable
      where PAYLOAD_MASS_KG_ = (select MAX(PAYLOAD_MASS_KG_) from spacetable)

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

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

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

# 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

```
%%sql select substr(date,4,2) as Month,Landing_Outcome,Booster_Version,launch_site from spacetable
      where substr(date,7,4) = '2015'
      and Landing_Outcome = 'Failure (drone ship)'
```

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

Done.

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

# 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 the grouped landing outcome in descending order.

```
%%sql SELECT landing_Outcome, COUNT(landing_Outcome) AS Count FROM spacetable  
      WHERE DATE(date) BETWEEN '2010-06-04' AND '2017-03-20'  
      GROUP BY landing_Outcome  
      ORDER BY Count DESC;
```

	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

A satellite view of Earth from space, showing the curvature of the planet and city lights at night. The background is a deep blue gradient.

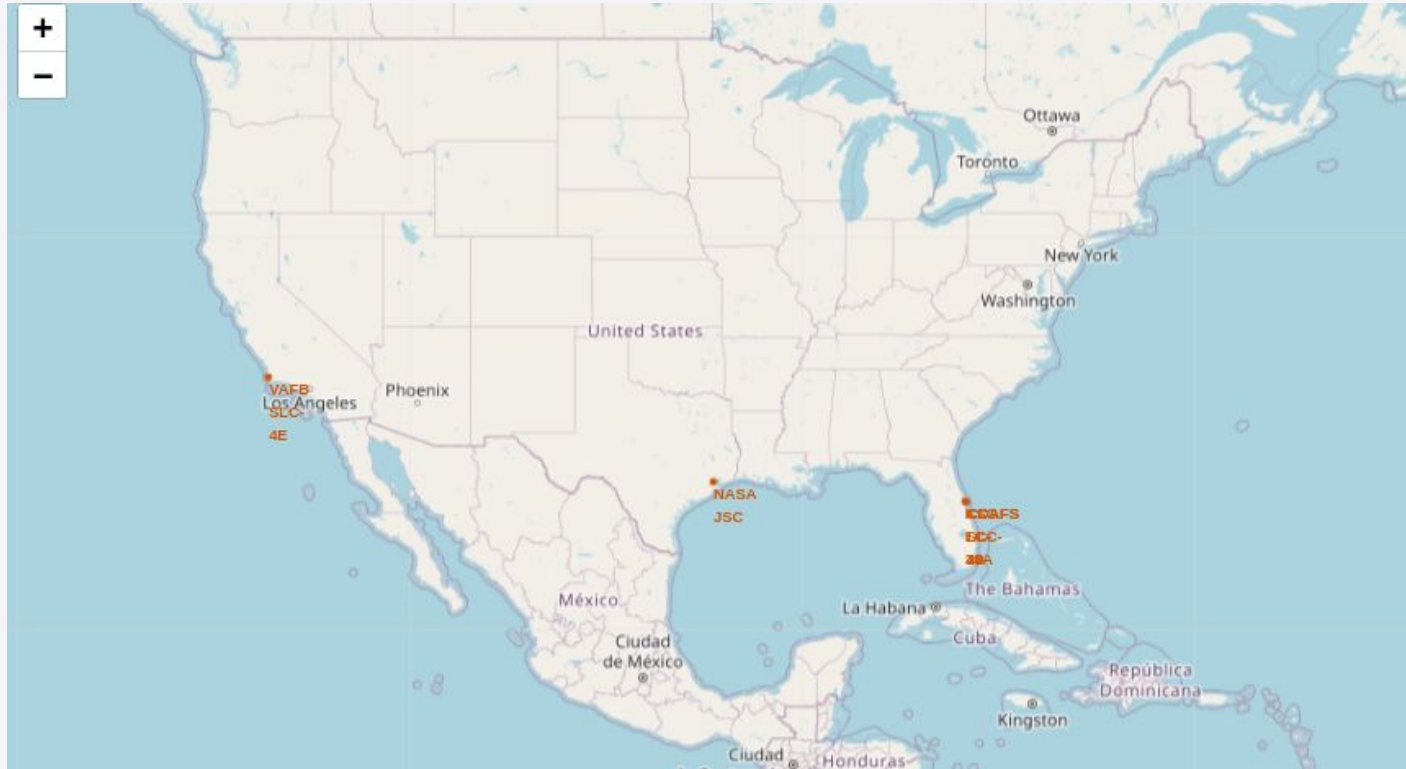
Section

3

# Launch Sites Proximities Analysis

# Marking all launch sites on a map

---

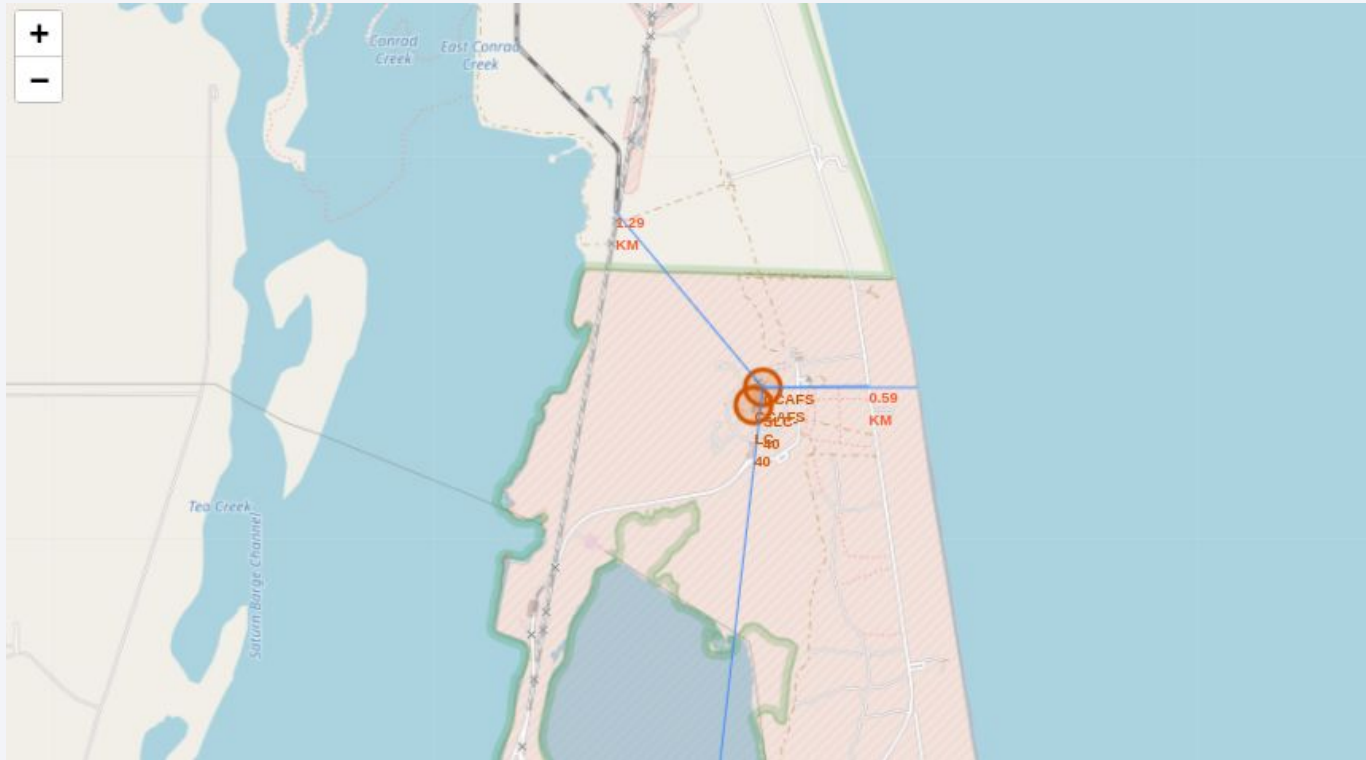






# The distances between a launch site to its proximities

---





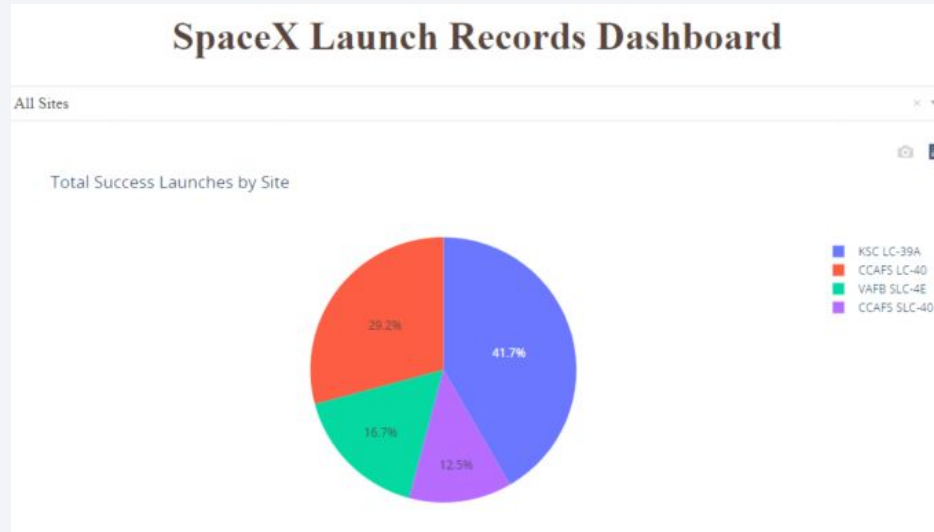
Section

4

# Build a Dashboard with Plotly Dash

# Launch success count for all sites

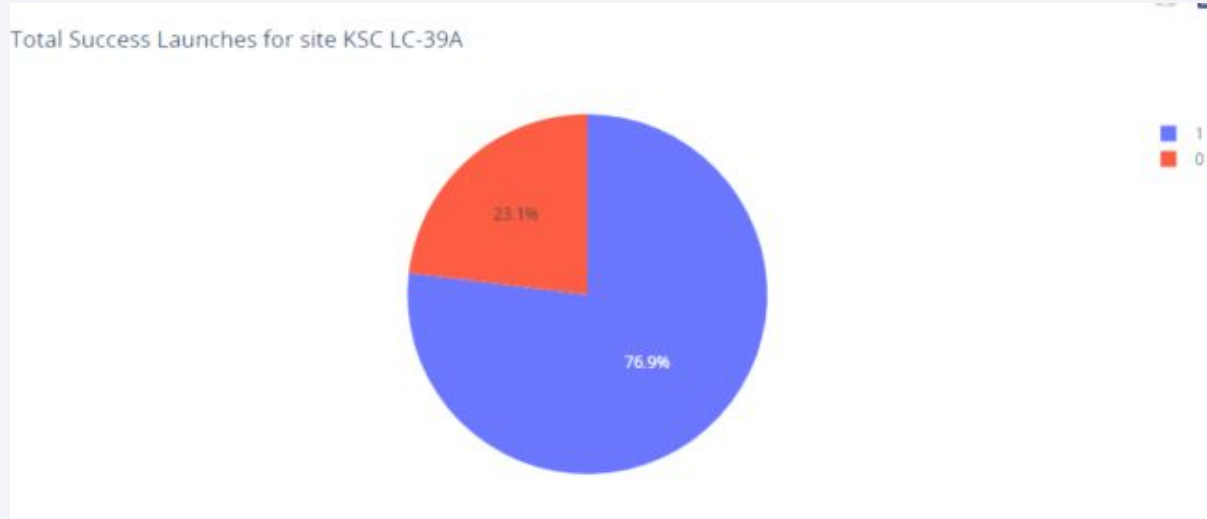
---



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

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

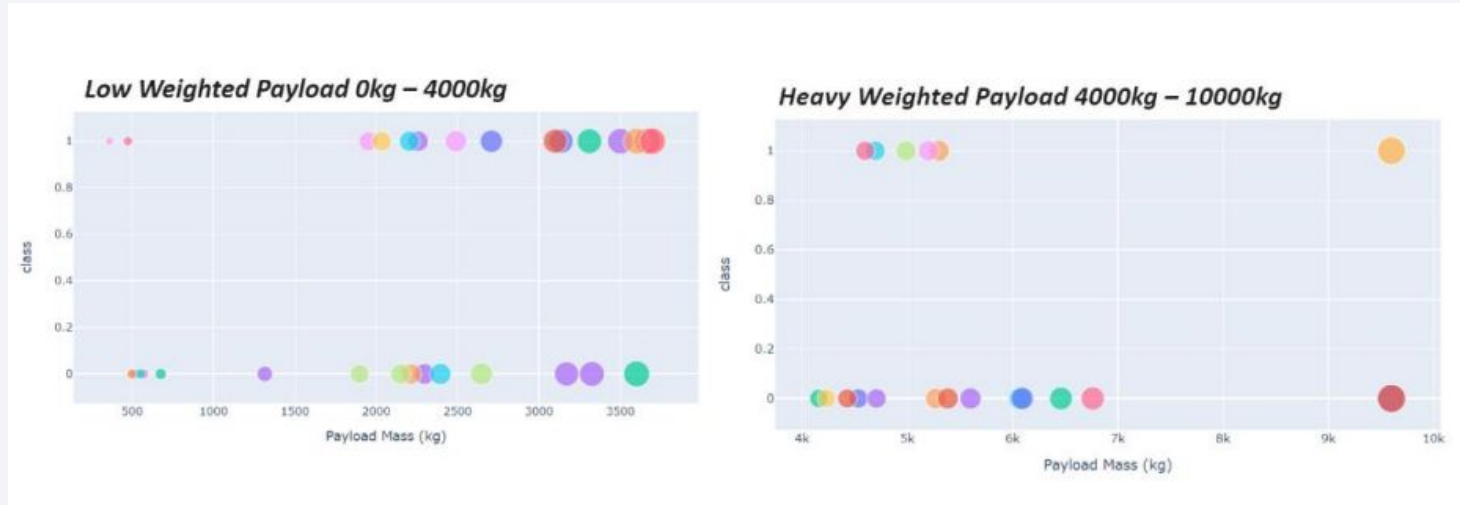
---



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

catter plot of Payload vs Launch Outcome for all sites, with different payload selected in the range slider

---



From these 2 plots, it can be shown that the success for massive payloads is lower than that for low payloads.

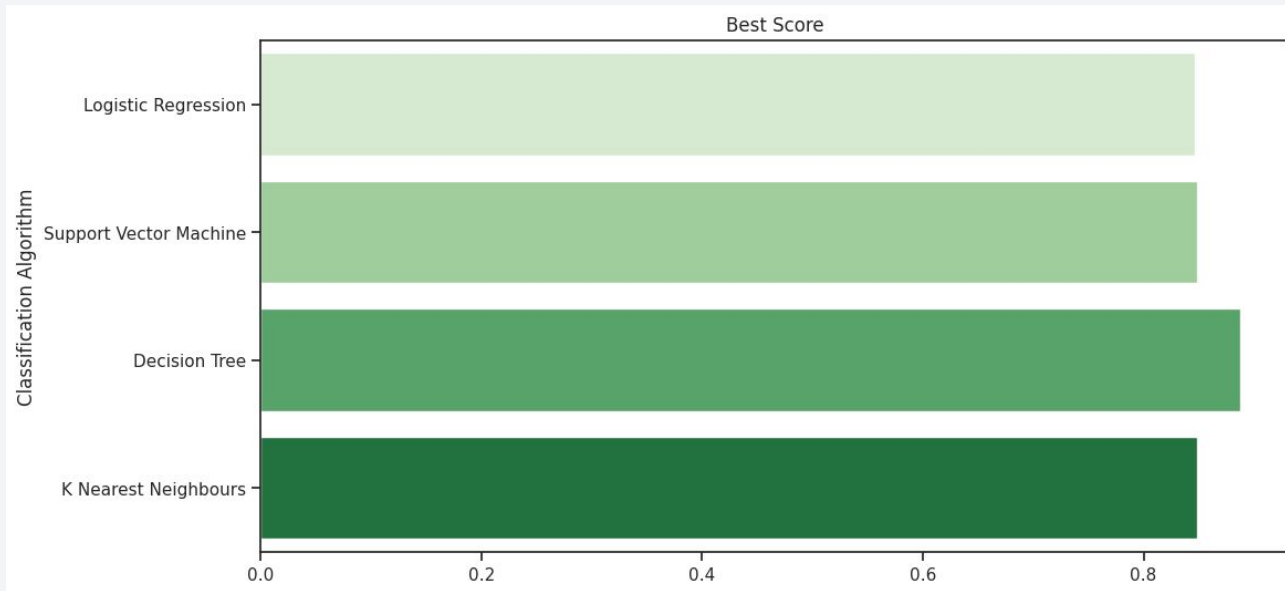


Section

5

# Predictive Analysis (Classification)

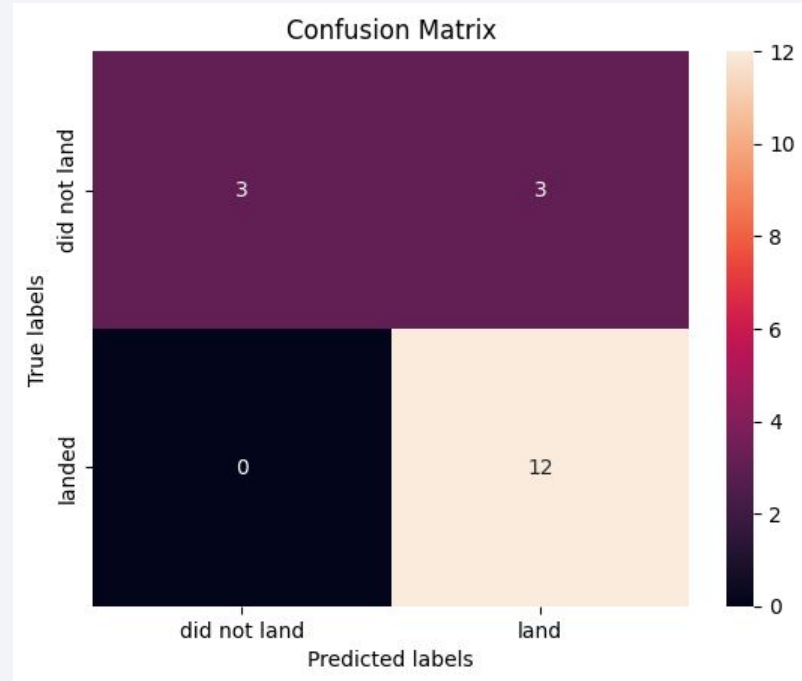
# Classification Accuracy



- 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 Best Score is 83.34%
  - The Accuracy Score is 88.75%

# 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!

