

# Winning Space Race with Data Science

Rufat Efendihev  
May 4, 2023



# Outline

---

- Executive Summary
- Introduction
- Methodology
- Results
- Conclusion



# Executive Summary

---

By using Python 3 and its libraries, historical data of **Falcon 9 launches** was retrieved through calls to the SpaceX's API and web scraping to the Wikipedia's article on List of Falcon 9 and Falcon Heavy launches.

Data was wrangled and explored through SQL queries, visualization and dashboarding, allowing to identify that SpaceX has a **success rate** of about **80%**; **KSC LC-39A** is the launch site with the highest number of successful missions, whereas **VAFB SLC-4E** is the one with the lowest; **ES-L1, GEO, HEO** and **SSO** were the orbit types with the highest success rates (**100%**) but with a small number of missions; **VLEO** was the most common orbit type in recent years; the **payload range** with the highest launch success rate was from **1952 kg** to **5300 kg**; and **FT** and **B5** are the **booster versions** with the highest launch success rate.

## Executive Summary

---

Finally, several classification models using Logistic Regression, Support Vector Machines, Decision Trees and K-Nearest Neighbors were set up and tuned, allowing to conclude that most of the launches from SpaceX will land successfully, thus suggesting that the cost of a Falcon 9 rocket launch should be set in 62 million dollars.



# Introduction

---

SpaceX offers Falcon 9 rocket launches with a cost of 62 million dollars while other providers cost upward of 165 million dollars each. A significant amount of the savings is due to SpaceX's ability to reuse the first stage.

In this context, the main goal of the present study is to predict the first stage landing of the SpaceX Falcon 9 rocket launch by using a classification model in order to determine the cost of a launch.

Furthermore, it is also desirable to obtain other insights from the SpaceX Falcon 9 rocket launches.



Section 1

# Methodology

# Methodology

---

## Executive Summary

Python 3 and its libraries were used in the entire project.

First, data was collected through REST calls to the SpaceX API using Requests and by performing web scraping to the Wikipedia's article on List of Falcon 9 and Falcon Heavy launches using BeautifulSoup.

Then, a process of data wrangling was performed using Pandas and Numpy in order to filter the data to the Falcon 9 launches only, imputing missing values with the mean, creating the landing outcome labels; as well as selecting the features for modeling and obtaining dummy variables for the categorical variables.

Next, an exploratory data analysis (EDA) was performed by means of visualization using Matplotlib and SQL.

# Methodology

---

## Executive Summary

After that, an **interactive visual analysis** using Folium and Plotly Dash was carried out.

Finally, a **predictive analysis** was performed by means of a Classification Model using Scipy and Scikit-learn. To do so:

- Data was split in a training and test sets.
- The model was built by using several machine learning techniques: Logistic Regression, Support Vector Machines (SVM), Decision Trees and K-Nearest Neighbors (KNN).
- The models were tuned by using GridSearchCV, in which several parameters were tested, and the data was cross-validated in a 10-fold scheme.
- The best model was selected based on the criteria of outcome of the confusion matrix, precision, recall, f1-score and accuracy.

# Data Collection

---

Data sets were collected through SpaceX API and Web Scraping in Python 3 by using the libraries Requests, Pandas, Numpy, BeautifulSoup, among others.



**Figure 1.** Data sources.

# Data Collection – SpaceX API

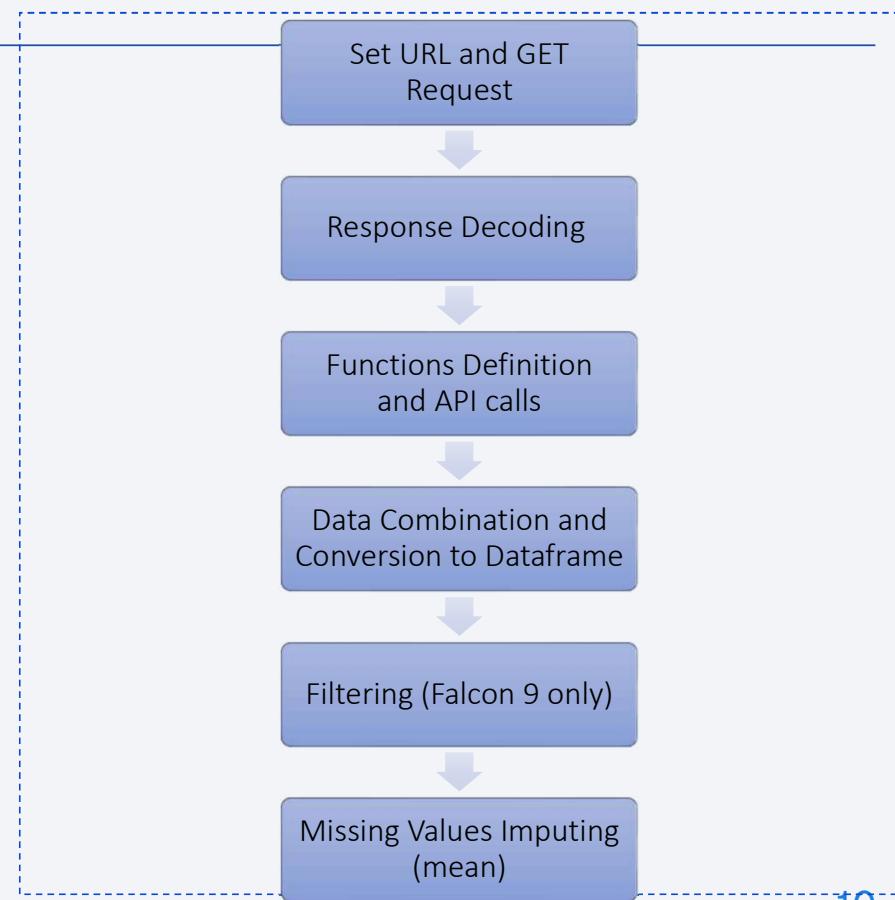
The purpose of this step was to make a get request to the SpaceX API and clean the requested data.

**Link of the data source:**

[https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBM-DS0321EN-SkillsNetwork/datasets/API\\_call\\_spacex\\_api.json](https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBM-DS0321EN-SkillsNetwork/datasets/API_call_spacex_api.json)

**GitHub URL:**

<https://github.com/erufat/IBM-SpaceX/blob/main/jupyter-labs-spacex-data-collection-api.ipynb>



**Figure 2.** Data collection process for SpaceX API.

# Data Collection - Scraping

The purpose of this step was to extract the Falcon 9 launch records from Wikipedia.

**Link of the data source:**

[https://en.wikipedia.org/wiki/List  
of Falcon\ 9\ and Falcon Heavy  
launches](https://en.wikipedia.org/wiki/List_of_Falcon_9_and_Falcon_Heavy_launches)

**GitHub URL:**

[https://github.com/erufat/IBM-  
SpaceX/blob/main/jupyter-labs-  
webscraping.ipynb](https://github.com/erufat/IBM-SpaceX/blob/main/jupyter-labs-webscraping.ipynb)



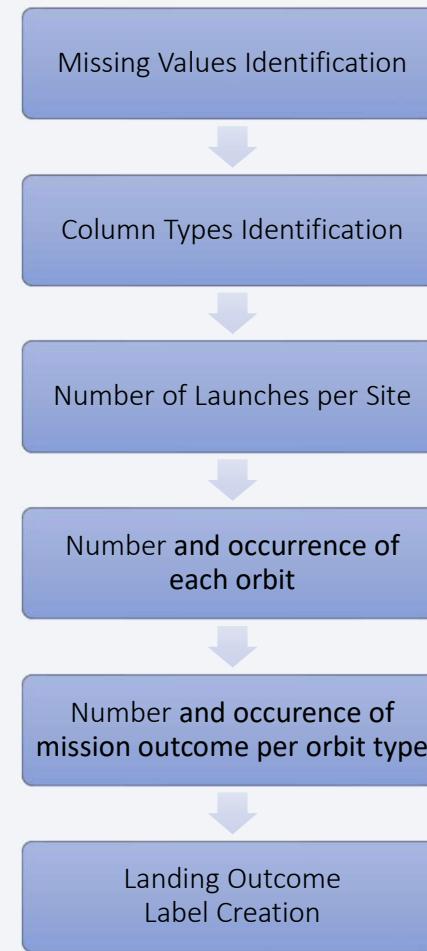
**Figure 3.** Data collection process for web scraping. 11

# Data Wrangling

The purpose of this step was to find some patterns in the data and determine the label for training supervised models.

**GitHub URL:**

[https://github.com/erufat/IBM-SpaceX/blob/main/IBM-DS0321EN-SkillsNetwork\\_labs\\_module\\_1\\_L3\\_labs-jupyter-spacex-data\\_wrangling\\_jupyterlite.jupyterlite.ipynb](https://github.com/erufat/IBM-SpaceX/blob/main/IBM-DS0321EN-SkillsNetwork_labs_module_1_L3_labs-jupyter-spacex-data_wrangling_jupyterlite.jupyterlite.ipynb)



**Figure 4.** Data wrangling process.

# EDA with Data Visualization

---

The purpose of this step was to perform Exploratory Data Analysis with Pandas and Matplotlib.

The following charts were plotted:

- **Scatter Plot of "Flight Number" vs. "Payload Mass":** To assess the likelihood of the first stage return in function of the Payload Mass and the Flight Number.
- **Scatter Plot of "Flight Number" vs. "Launch Site":** To assess the relationship among Launch Site, Flight Number, and the first stage return.
- **Scatter Plot of "Payload Mass" vs. "Launch Site":** To assess the relationship among Launch Site, Payload Mass and the first stage return.

# EDA with Data Visualization

---

- **Bar chart of the "Orbit" "Success Rate":** To assess the sucess rate of each orbit.
- **Scatter Plot of "Flight Number" vs. "Orbit":** To assess the relationship among Flight Number, Orbit, and the first stage return.
- **Scatter Plot of "Payload Mass" vs. "Orbit":** To assess the relationship among Payload Mass, Orbit, and the first stage return.
- **Line Chart of "Payload Mass" vs. "Orbit":** To visualize the launch success yearly trend.

**GitHub URL:** [https://github.com/erufat/IBM-SpaceX/blob/main/IBM-DS0321EN-SkillsNetwork\\_labs\\_module\\_2\\_jupyter-labs-eda-dataviz.ipynb.jupyterlite.ipynb](https://github.com/erufat/IBM-SpaceX/blob/main/IBM-DS0321EN-SkillsNetwork_labs_module_2_jupyter-labs-eda-dataviz.ipynb.jupyterlite.ipynb)

# Feature Engineering

---

In addition, Feature Engineering was carried out to select the features that were used in the classification model. The following ones were selected:

```
'FlightNumber', 'PayloadMass', 'Orbit', 'LaunchSite', 'Flights',  
'GridFins', 'Reused', 'Legs', 'LandingPad', 'Block', 'ReusedCount',  
'Serial'
```

Then, dummy variables were created for the categorical variables Orbit, LaunchSite, LandingPad, and Serial using the `get_dummies()` function.

Finally, the whole dataset was casted to float64 datatype.

# Build an Interactive Map with Folium

---

An Interactive Map with the Folium library was created, and the following maps objects were created and added to the map:

- **Markers and circles for all launch sites:** To visualize the location of the launch sites and assess their characteristics.
- **Marker Clusters for success/failed launches for each site:** To assess if there was a relationship among the launch sites and their success/fail missions.
- **Markers and lines between a launch site to its proximities:** To assess the characteristics of the proximities to the launch sites.
- **GitHub URL:** [https://github.com/erufat/IBM-SpaceX/blob/main/IBM-DS0321EN-SkillsNetwork\\_labs\\_module\\_3\\_lab\\_jupyter\\_launch\\_site\\_location.ipynb](https://github.com/erufat/IBM-SpaceX/blob/main/IBM-DS0321EN-SkillsNetwork_labs_module_3_lab_jupyter_launch_site_location.ipynb)

# Build a Dashboard with Plotly Dash

---

An Interactive Dashboard with the Plotly and Dash libraries was built, and the following plots/graphs and interactions were added to the dashboard:

- **Dropdown:** In order to enable the Launch Site selection for the user.
- **Pie chart:** To show the Total Success Launches by Site.
- **Range Slider:** To enable the Payload Mass selection.
- **Scatter Plot:** To show the relationship among Payload Mass and mission outcomes for the selected Sites.

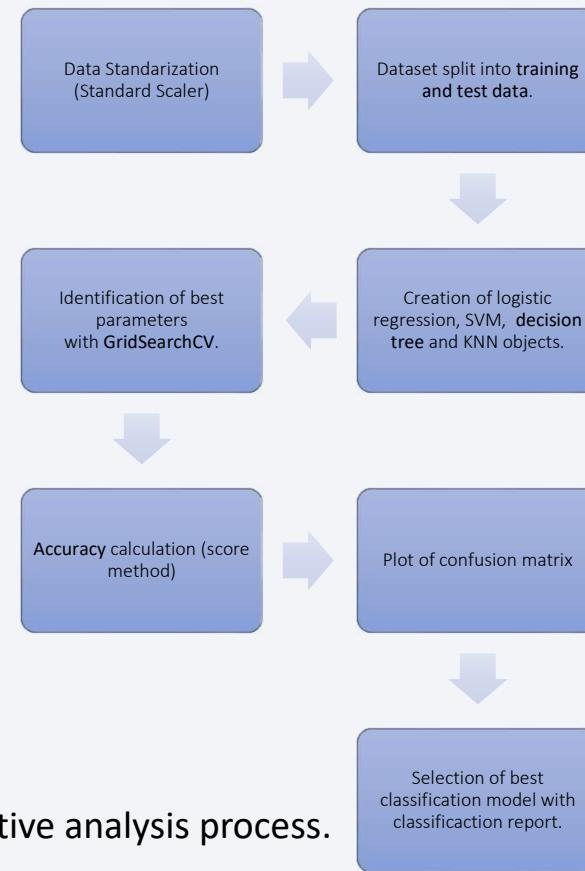
GitHub URL: [https://github.com/erufat/IBM-SpaceX/blob/main/SpaceX\\_Dashboard.ipynb](https://github.com/erufat/IBM-SpaceX/blob/main/SpaceX_Dashboard.ipynb)

# Predictive Analysis (Classification)

The purpose of this step was to build a Classification model for predicting the landing of the first phase; as well as finding the best hyperparameters and the best model.

**GitHub URL:**

[https://github.com/erufat/IBM-SpaceX/blob/main/IBM-DS0321EN-SkillsNetwork labs module 4 SpaceX Machine Learning Prediction Part 5.jupyterlite.ipynb](https://github.com/erufat/IBM-SpaceX/blob/main/IBM-DS0321EN-SkillsNetwork%20labs%20module%204%20SpaceX%20Machine%20Learning%20Prediction%20Part%205.jupyterlite.ipynb)



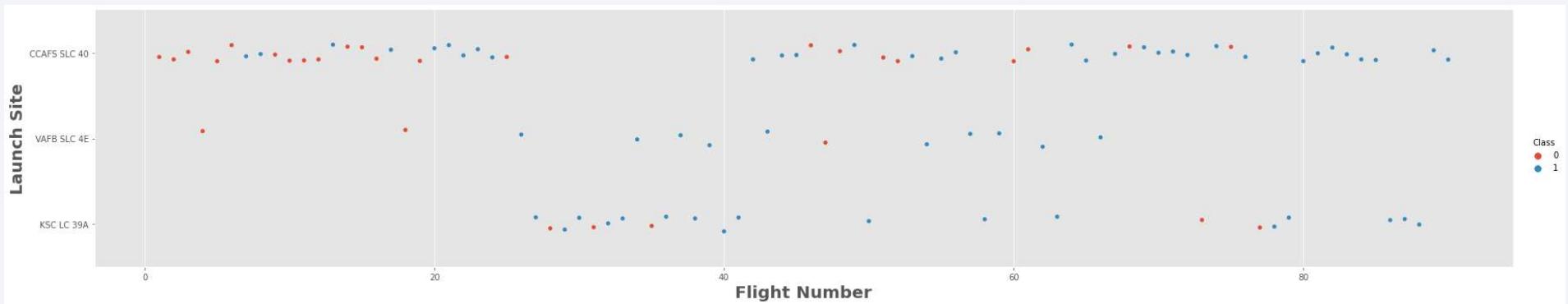
**Figure 5.** Predictive analysis process.

The background of the slide features a dynamic, abstract pattern of glowing lines. These lines are primarily blue and red, with some green and purple highlights. They appear to be moving in a three-dimensional space, creating a sense of depth and motion. The lines are thick and have a slight glow, making them stand out against the dark background.

Section 2

## Insights drawn from EDA

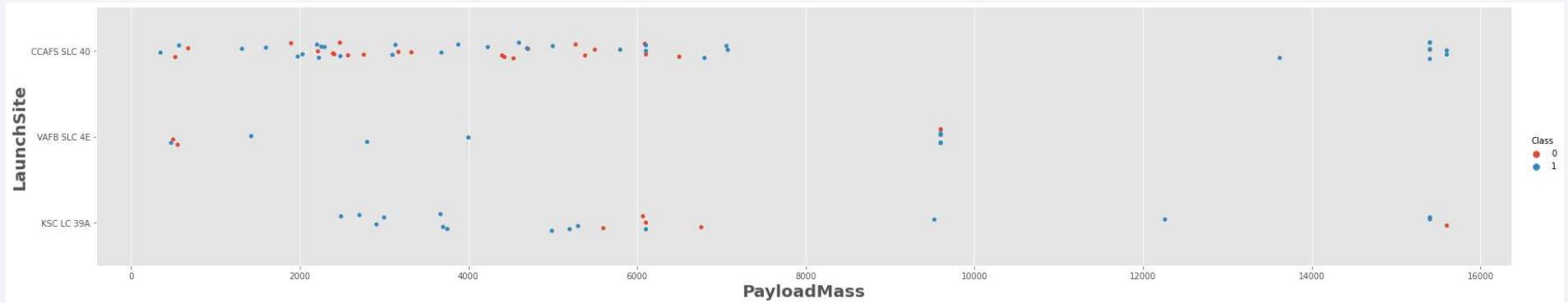
# Launch Site vs. Flight Number



**Figure 6.** Launch Site vs. Flight Number.

- CCAFS SLC 40 is the most used Launch Site and the VAFB SLC-4E is the least used one.
- The success rate has improved over the years as represented by the Flight number.
- KSC LC-39A is the launch site with the highest number of successful missions.

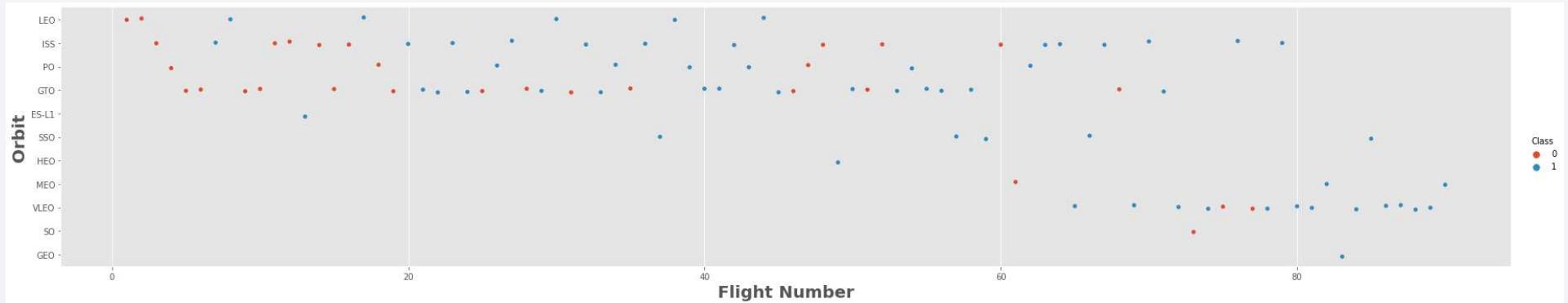
# Launch Site vs. Payload



**Figure 7.** Launch Site vs. Payload.

- CCAFS SLC 40 and KSC LC-39A are used for a broad range of payload masses, from light to extra heavy ones.
- VAFB SLC-4E is only used for light and medium payloads, which might explain why it is the less used launch site.
- The missions tend to be more successful with heavier payloads.

# Orbit Type vs. Flight Number



**Figure 8.** Orbit Type vs. Flight Number.

- LEO, ISS, PO, GTO and VLEO are the most common orbit types.
- EO, PO, GTO and VLEO are less common in recent flight numbers.
- VLEO is the most common orbit type in recent flight numbers.
- Earlier flight numbers show a tendency to fail regardless of the orbit type.

# Orbit Type vs. Payload

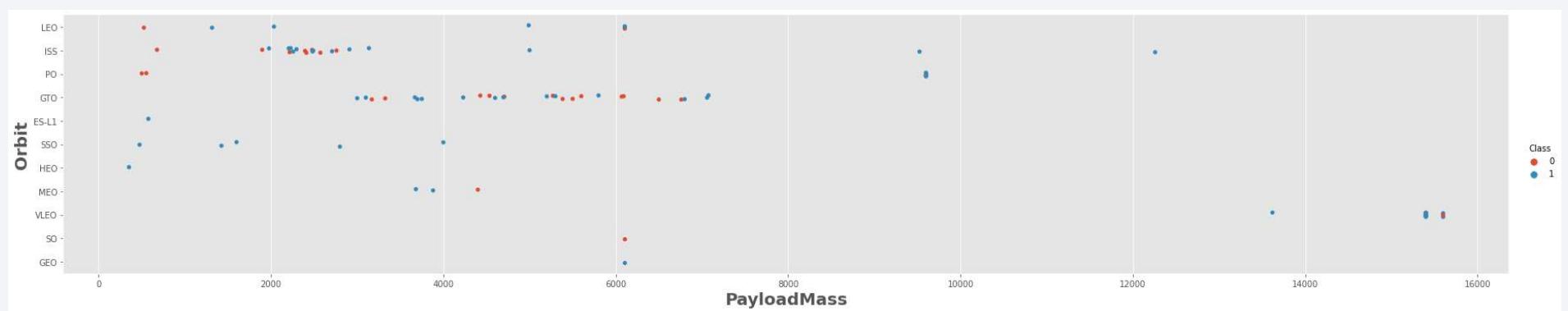


Figure 9. Orbit Type vs. Payload.

- Heavier payloads are selected for VLEO, IS and PO orbit types and have a high rate of success.
- It appears that ES-L1, SSO and HEO are used with light payloads (from 0 to 4000 kg) and have a high rate of success.
- The rest of the orbit types are used with light and medium payloads (from 0 to 7000 kg).

# All Launch Site Names

---

CCAFS SLC-40 is the most used Launch Site and the VAFB SLC-4E is the least used.

```
In [ ]: %%sql
```

```
SELECT Launch_Site, COUNT(Launch_Site) AS Count FROM SPACEXTBL GROUP BY Launch_Site;
```

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

```
Done.
```

```
Out[ ]: Launch_Site  Count
```

Launch_Site	Count
CCAFS LC-40	26
CCAFS SLC-40	34
KSC LC-39A	25
VAFB SLC-4E	16

**Figure 10.** SQL Query with Launch Site Count.

# Launch Site Names Begin with 'CCA'

The launch site whose name begins with the string 'CCA' is CCAFS SLC-40.

2010-04-06	18:45:00	F9 v1.0 B0003	CCAFS LC-40	Dragon Spacecraft Qualification Unit	0	LEO	SpaceX	Success	Failure (parachute)
2010-08-12	15:43:00	F9 v1.0 B0004	CCAFS LC-40	Dragon demo flight C1, two CubeSats, barrel of Brouere cheese	0	LEO (ISS)	NASA (COTS) NRO	Success	Failure (parachute)
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
2012-08-10	00:35:00	F9 v1.0 B0006	CCAFS LC-40	SpaceX CRS-1	500	LEO (ISS)	NASA (CRS)	Success	No attempt
2013-01-03	15:10:00	F9 v1.0 B0007	CCAFS LC-40	SpaceX CRS-2	677	LEO (ISS)	NASA (CRS)	Success	No attempt

**Figure 11.** SQL Query with Launch Site whose name begins with 'CCA'.

# Total Payload Mass

---

The total payload carried by boosters from NASA is 45 596 kg.

```
In [ ]: %%sql  
  
SELECT SUM(PAYLOAD_MASS__KG_) FROM SPACEXTBL WHERE Customer = 'NASA (CRS)';  
* sqlite:///my_data1.db  
Done.  
Out[ ]: SUM(PAYLOAD_MASS_KG_)  
_____  
45596
```

**Figure 12.** SQL Query with Total Payload Carried by Booster from NASA.

# Average Payload Mass by F9 v1.1

---

The average payload mass carried by booster version F9 v1.1 is 2534.66 kg, which is in the range of the light payloads.

```
In [ ]: %%sql
SELECT AVG(PAYLOAD_MASS__KG_) FROM SPACEXTBL WHERE Booster_Version LIKE 'F9 v1.1%';
* sqlite:///my_data1.db
Done.

Out[ ]: AVG(PAYLOAD_MASS__KG_)
2534.666666666665
```

**Figure 13.** SQL Query with average payload mass carried by booster version F9 v1.1.

# First Successful Ground Landing Date

---

The first successful landing outcome on ground pad was achieved on 2015-12-22.

```
In [ ]: %%sql
SELECT MIN(Date) FROM SPACEXTBL WHERE "Landing _Outcome" = 'Success (ground pad)';
* sqlite:///my_data1.db
Done.

Out[ ]: MIN(Date)
2015-12-22
```

**Figure 14.** SQL Query with first successful landing outcome on ground pad.

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

---

- The names of boosters which have successfully landed on drone ship and had payload mass greater than 4000 but less than 6000 are F9 FT B1022, F9 FT B1026, F9 FT B1021.2 and F9 FT B1031.2.

```
In [ ]: %%sql
SELECT Booster_Version FROM SPACEXTBL WHERE "Landing _Outcome" = 'Success (drone ship)' AND PAYLOAD_MASS__KG_ > 4000 AND PAYLOAD_MASS__KG_ < 6000;
* sqlite:///my_data1.db
Done.

Out[ ]: Booster_Version
F9 FT B1022
F9 FT B1026
F9 FT B1021.2
F9 FT B1031.2
```

**Figure 15.** SQL Query with names of boosters which have successfully landed on drone ship and had payload mass greater than 4000 but less than 6000.

## Total Number of Successful and Failure Mission Outcomes

---

The total number of successful and failure mission outcomes are 100 and 1, respectively.

```
In [ ]: %%sql
SELECT Mission_Outcome, COUNT(Payload) FROM SPACEXTBL GROUP BY Mission_Outcome;
* sqlite:///my_data1.db
Done.

Out[ ]:   Mission_Outcome  COUNT(Payload)
          Failure (in flight)      1
          Success                  98
          Success                  1
          Success (payload status unclear)  1
```

**Figure 16.** SQL Query with total number of successful and failure mission outcomes.

# Boosters Carried Maximum Payload

---

The names of the booster which have carried the maximum payload mass are shown in the image below:

List the names of the booster\_versions which have carried the maximum payload mass. Use a subquery

```
In [ ]: %%sql
SELECT Booster_Version FROM SPACEXTBL WHERE PAYLOAD_MASS__KG_ = (SELECT MAX(PAYLOAD_MASS__KG_) FROM SPACEXTBL);
* sqlite:///my_data1.db
Done.
```

Out[ ]: 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

**Figure 17.** SQL Query with names of the booster which have carried the maximum payload mass.

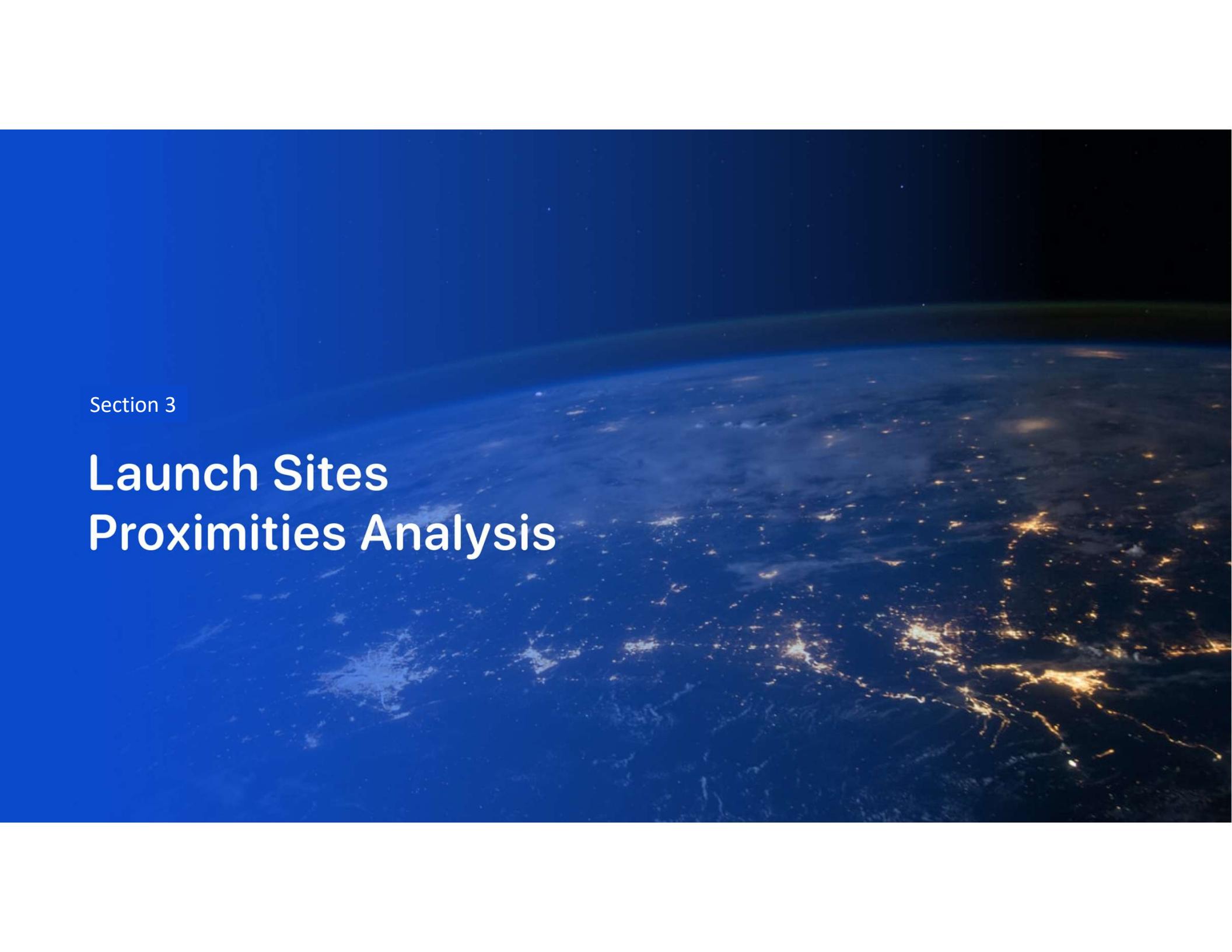
## Rank Landing Outcomes Between 2010-06-04 and 2017-03-20

The rank of 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 is shown in the image below:

```
In [ ]: %%sql
SELECT SUBSTR(Date,1,4) AS Year, COUNT("Landing _Outcome") AS Successful_landing_outcomes
FROM SPACEXTBL
WHERE DATE > '2010-06-04' AND DATE < '2017-03-20' AND "Landing _Outcome" LIKE 'Success%'
GROUP BY SUBSTR(Date,1,4)
ORDER BY SUBSTR(Date,1,4) DESC;
* sqlite:///my_data1.db
Done.

Out[ ]: Year  Successful_landing_outcomes
2017              4
2016              5
2015              1
```

**Figure 18.** SQL Query with rank of count of landing outcomes between the date 2010-06-04 and 2017-03-20.

The background of the slide is a photograph taken from space at night. It shows the curvature of the Earth's horizon against the dark void of space. City lights are visible as numerous small white and yellow dots, primarily concentrated in the lower right quadrant where major urban centers like North America are located. In the upper left quadrant, the green and blue glow of the aurora borealis (Northern Lights) is visible, appearing as a horizontal band of light.

Section 3

# Launch Sites Proximities Analysis

# Launch Sites Localization

Launch sites are located on the coasts and close to the Ecuador.

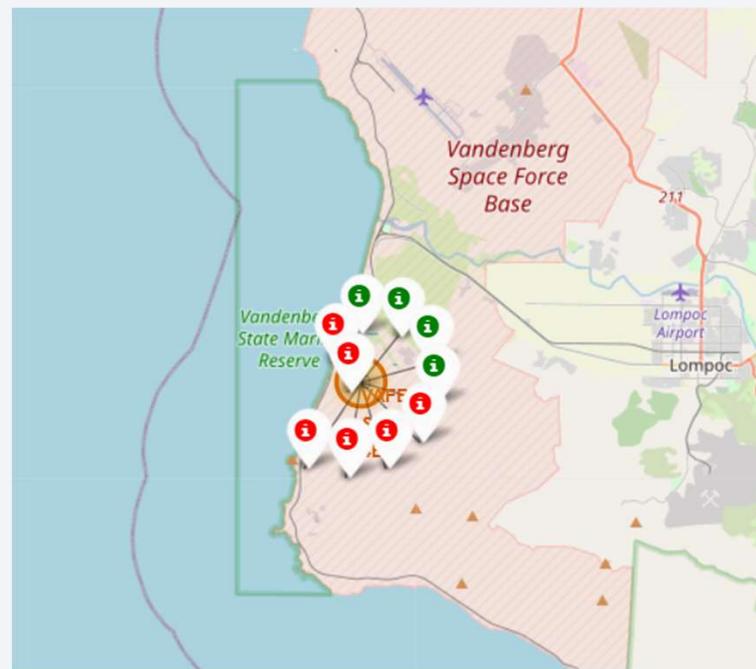


Figure 19. Map with the localization of launch sites.

## Launch Outcomes Per Launch Site

---

The Vandenberg (VAFB SLC-4E) Launch Site is the least used and its number of successful landings is fairly low.

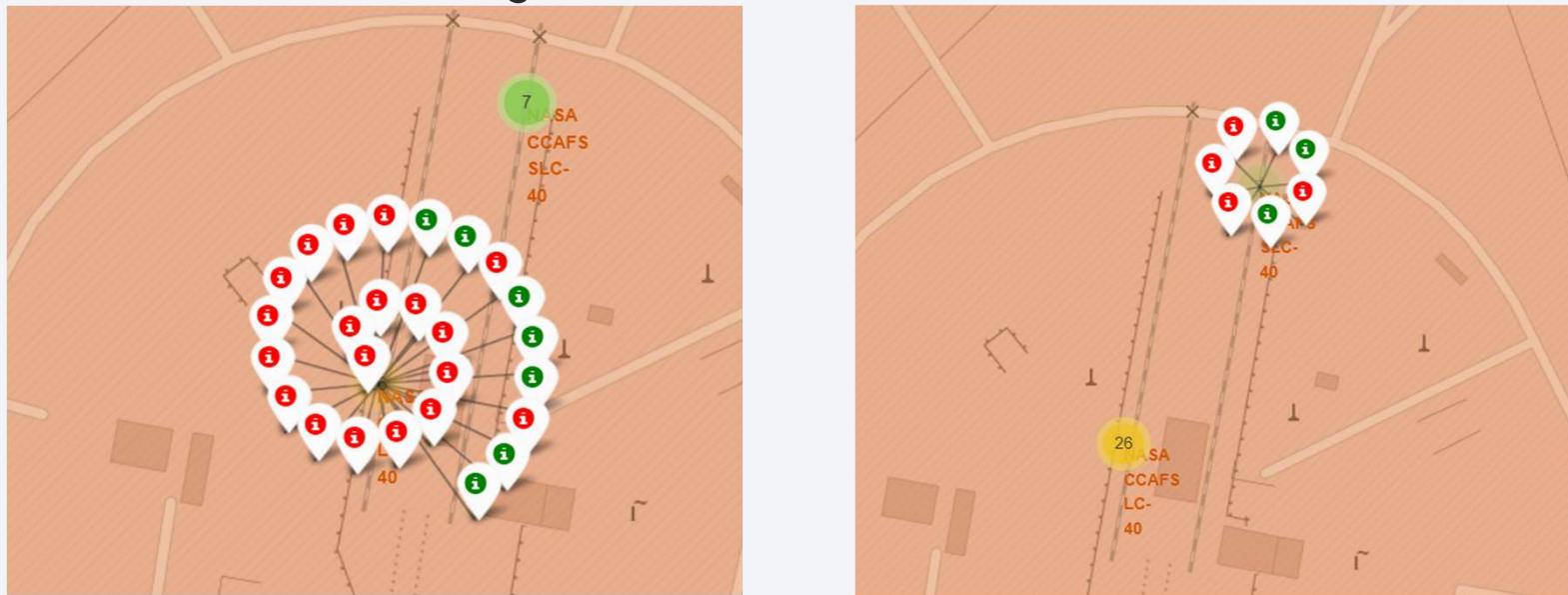


**Figure 20.** Map with the localization of VAFB SLC-4E launch site and marks of successful and unsuccessful landings.

## Launch Outcomes Per Launch Site

---

The Cape Canaveral (CCAFS LC-40 and CCAFS SLC-40) is the most used Launch Site, even though its number of failed landings is slightly superior to the its number of successful landings.

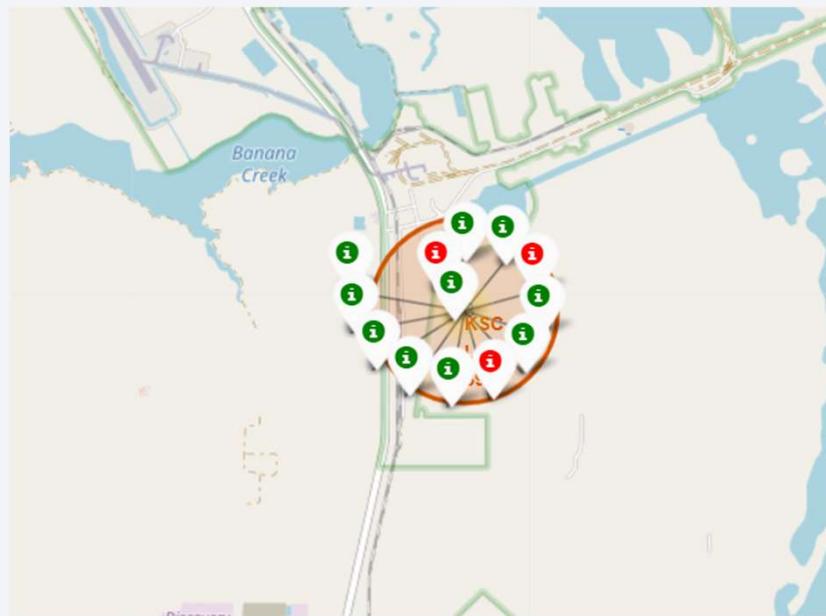


**Figure 21.** Map with the localization of CCAFS LC-40 and CCAFS SLC-40 launch site and marks of successful and unsuccessful landings.

# Launch Outcomes Per Launch Site

---

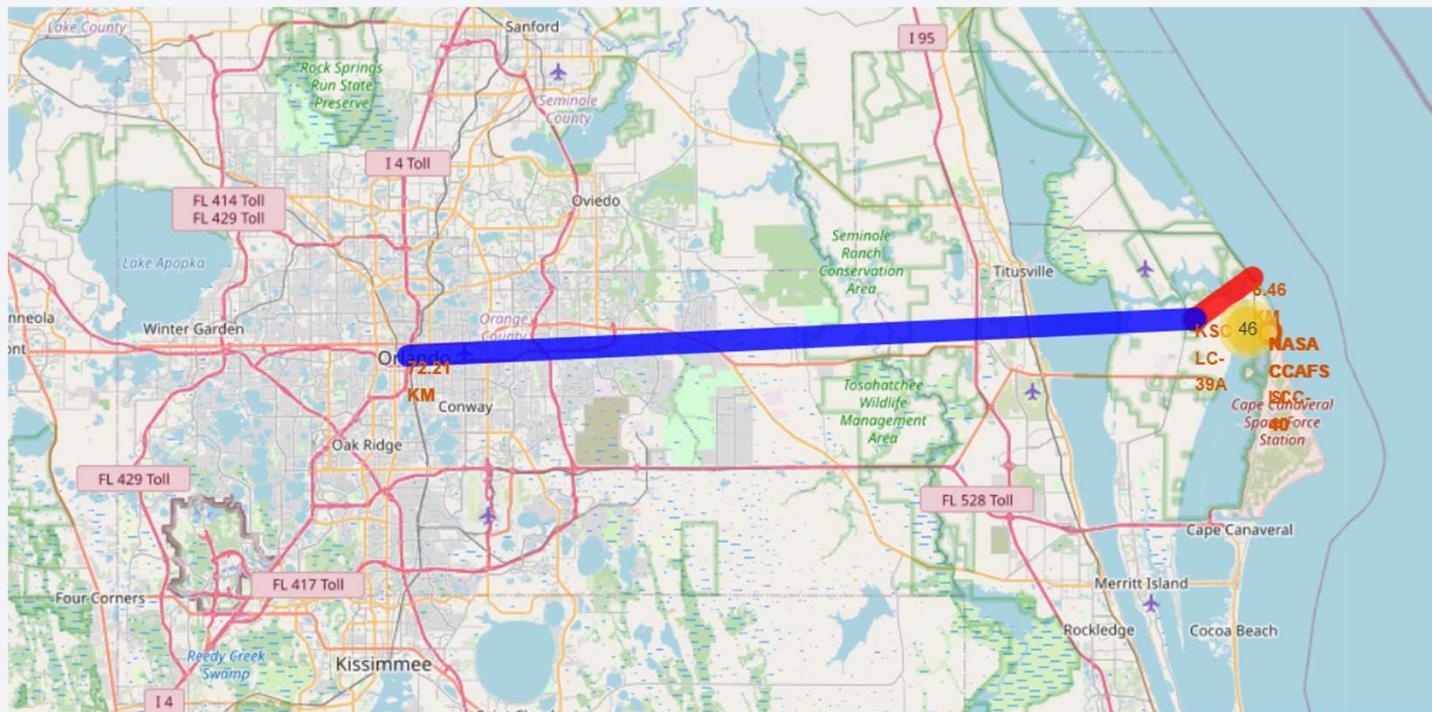
The Kennedy Space Center (**KSC LC-39A**) has the highest number of successful missions. Nonetheless, it is not the most used Launch Site.



**Figure 22.** Map with the localization of KSC LC-39A launch site and marks of successful and unsuccessful landings.

# KSC LC-39A distance to the sea and Orlando, FL.

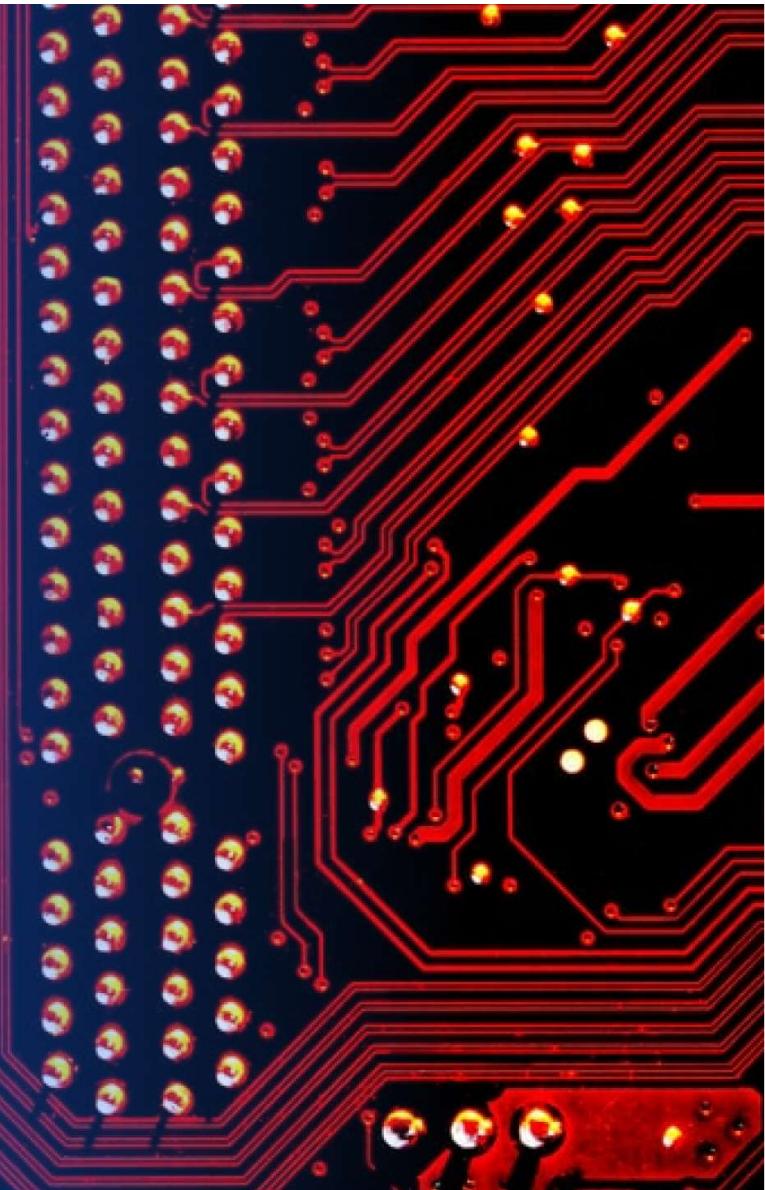
KSC LC-39A is 6.46 KM far from the Atlantic Sea, and 72.21 KM far from Orlando, FL.



**Figure 23.** Map with the localization of KSC LC-39A launch site and distances to the sea and Orlando, FL.

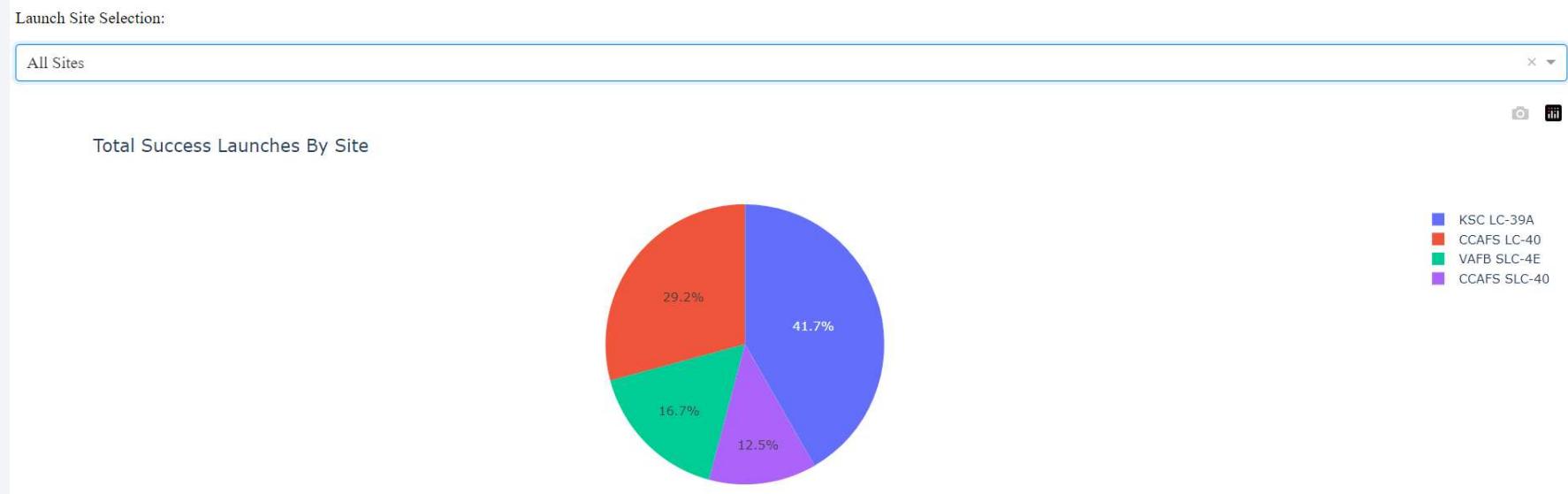
Section 4

# Build a Dashboard with Plotly Dash



## Launch Success Per Site

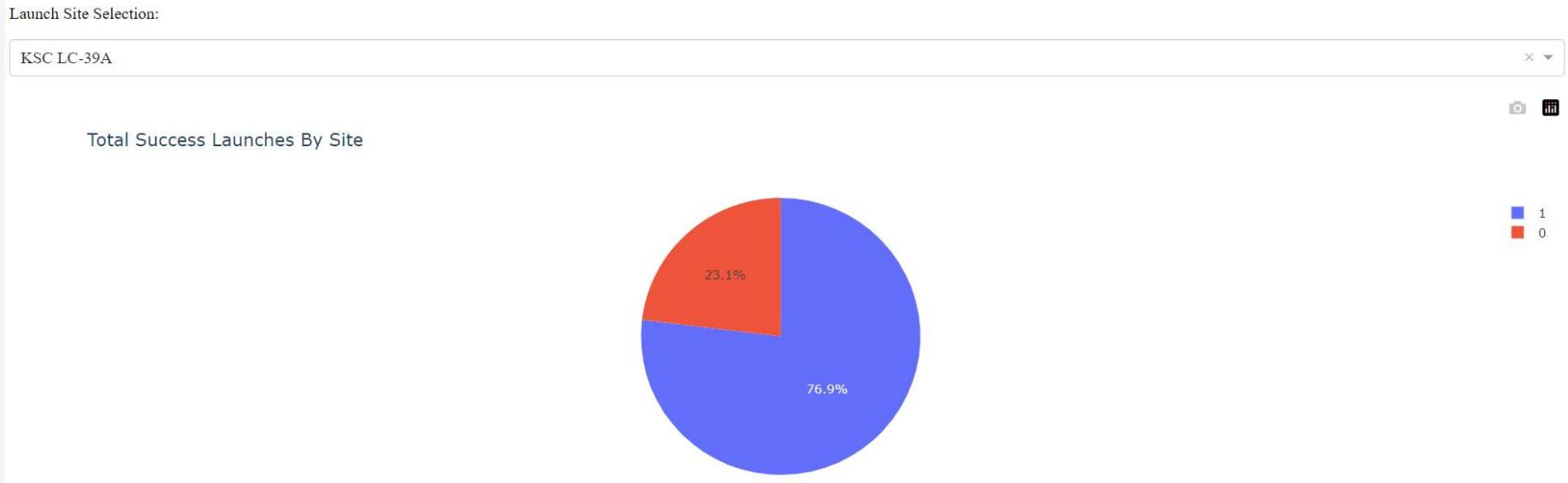
KSC LC-39A has the highest success launch rates, while VAFB SLC-4E holds the lowest. It is important to note that CCAFS LC-40 and CCAFS SLC-40 refer to the same facility.



**Figure 24.** Launch Success Per Site in Dashboard.

# Highest Launch Success Ratio

The Highest Launch Success Ratio is hold by the KSC LC-39A launch site, with a **76.9% of success** and only **23.1%** of failure.



**Figure 25.** Launch Success for KSC LC-39A in Dashboard.

# Payload vs. Launch Outcome

The Payload vs. Launch Outcome allowed to identify the payload ranges with the highest and the lowest success rates as well as the more successful F9 Booster versions (see next slides).

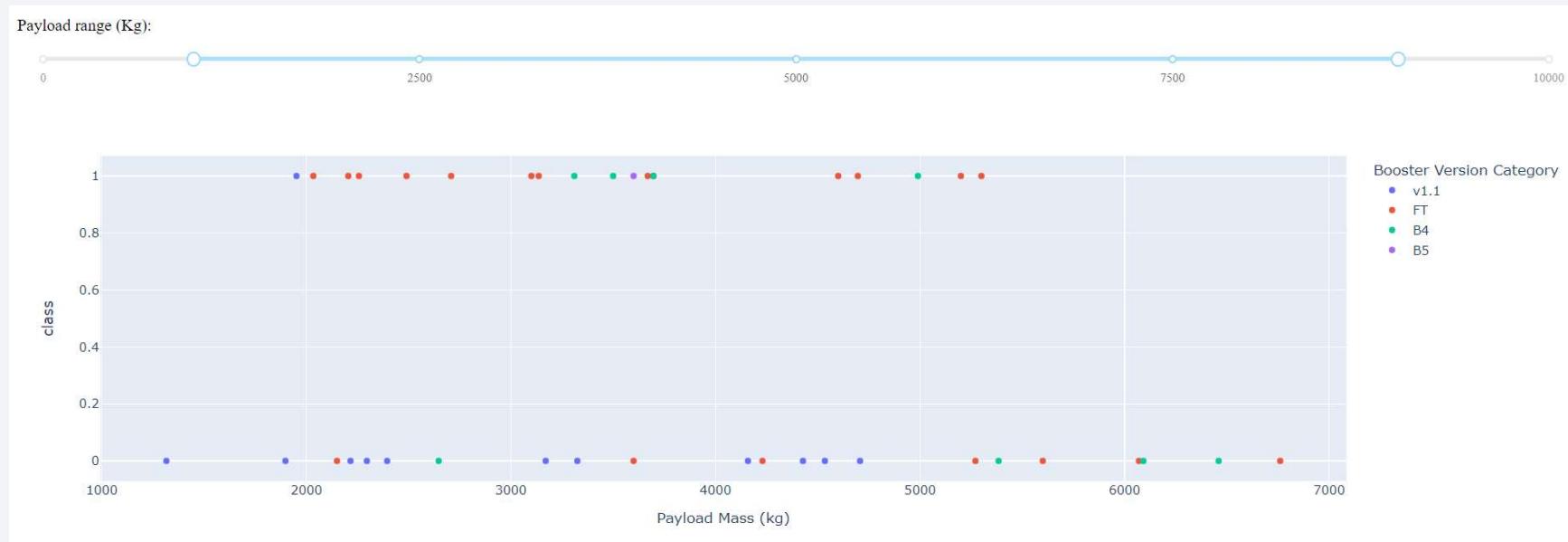


Figure 26. Dashboard's Scatter Plot of Payload Mass vs Class with Range Slider for Payload Mass 2.

# Payload vs. Launch Outcome

Which payload range(s) has the highest launch success rate? From 1952 kg to 5300 kg.



**Figure 27.** Payload range with the highest launch success rate.

# Payload vs. Launch Outcome

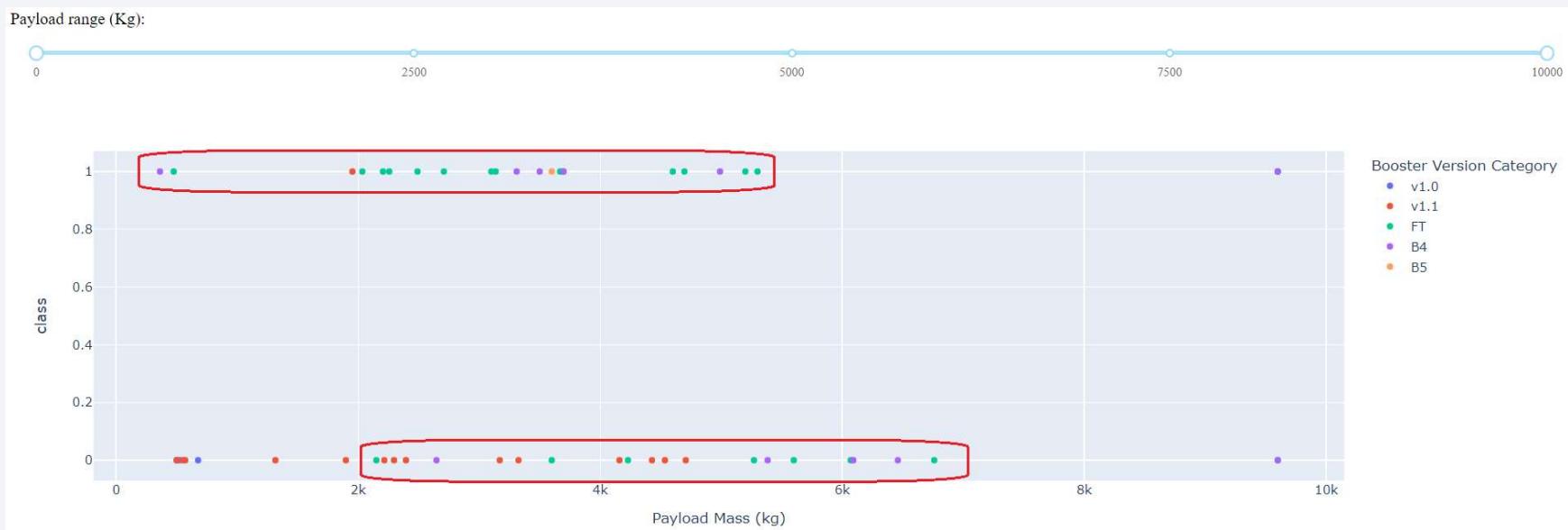
- Which payload range(s) has the lowest launch success rate? **Below 1952 kg and above 5300 kg.**



**Figure 28.** Payload ranges with the lowest launch success rate.

# Best Booster versions

Which F9 Booster version (v1.0, v1.1, FT, B4, B5, etc.) has the highest launch success rate? **The FT and the B5 Booster versions.**



**Figure 29.** F9 Booster versions with the highest launch success rate.

A blurred photograph of a tunnel, likely from a moving vehicle, showing motion streaks in shades of blue, white, and yellow. The perspective curves away from the viewer.

Section 5

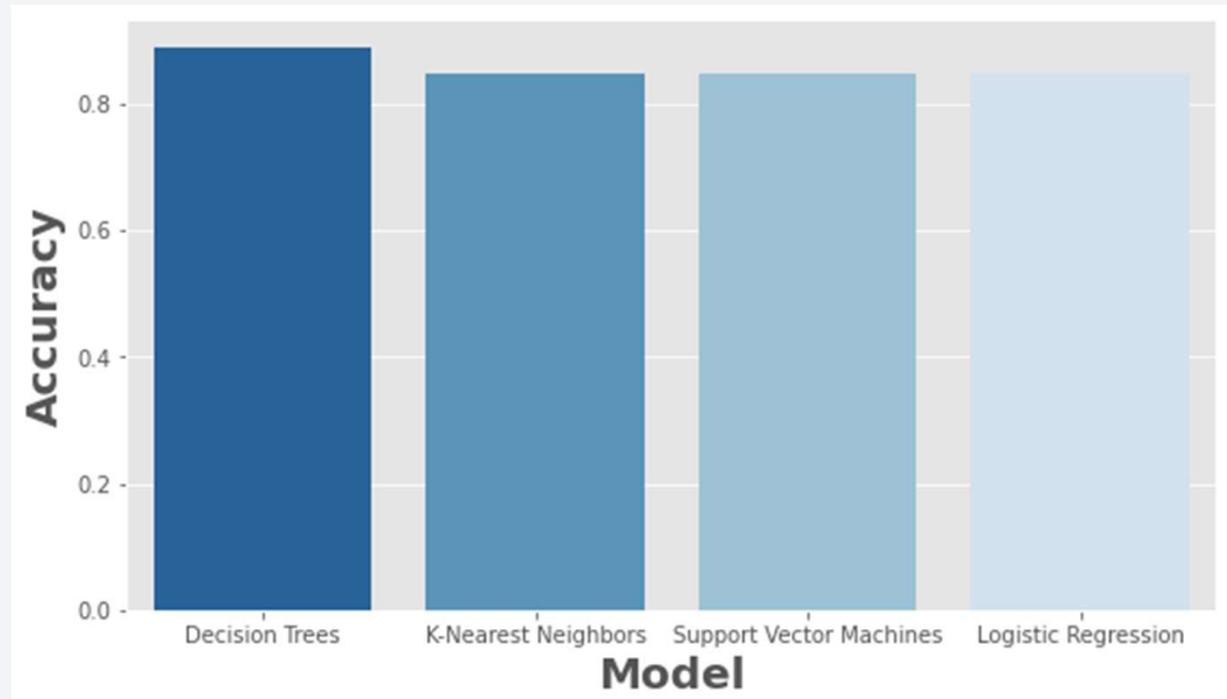
## Predictive Analysis (Classification)

# Classification Accuracy

---

The model with the highest classification accuracy was the one built with **Decision Trees** with an accuracy of **88.9%**, when fitted with the best parameters:

- {'criterion': 'gini',  
 'max\_depth': 2,  
 'max\_features': 'sqrt',  
 'min\_samples\_leaf': 1,  
 'min\_samples\_split': 2,  
 'splitter': 'best'}



**Figure 30.** Accuracy of the built classification models.

# Confusion Matrix

---

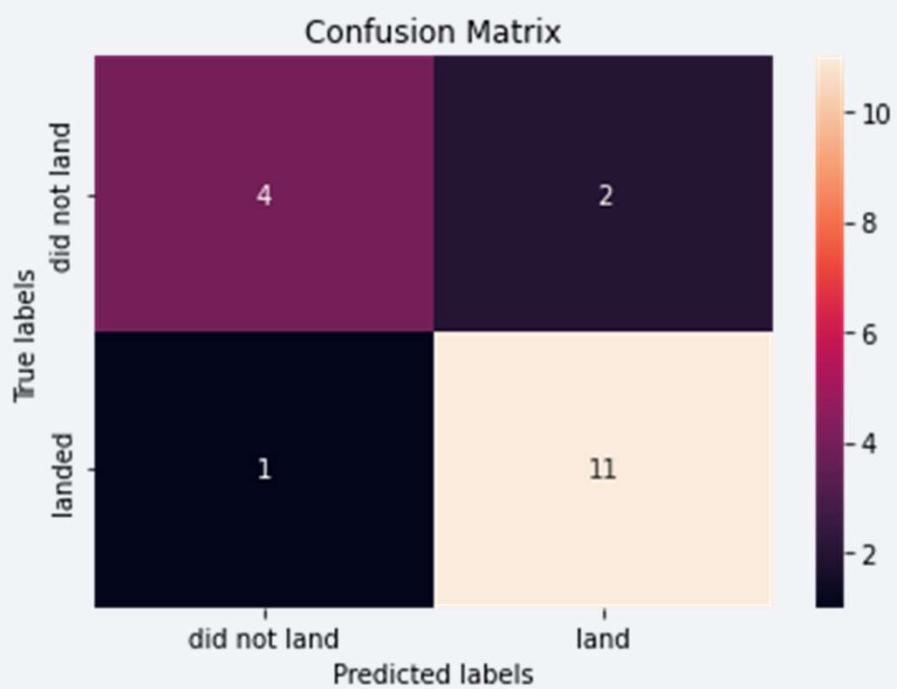


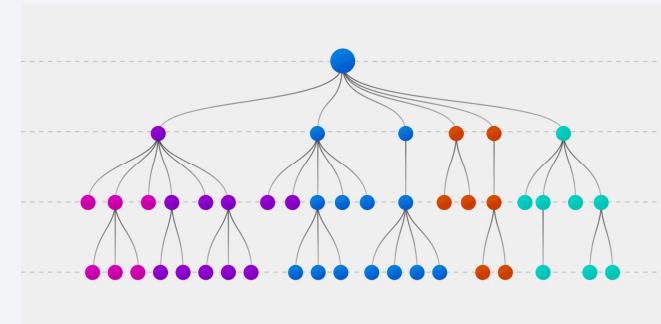
Figure 31. Decision trees confusion matrix.

- The confusion matrix suggests that the built model is better at predicting successful landings than unsuccessful landings.
- So, there is room for improvement in terms of reducing both false positives and false negatives.
- In spite of the above, this results suggest that, indeed, most of the launches from SpaceX will land successfully.

# Conclusions

---

- Firstly, a classification model using Decision Trees was developed with a precision of **88.9%** and whose outcome indicates that **most of the launches from SpaceX will land successfully**.
- On the other hand, analysis from the historical launching data suggests that SpaceX has gotten better at launching and its success rate has stabilized since 2017 in about **80%**.
- Thus, the present study suggests that the cost of the Falcon 9 rocket launches should be set at **62 million dollars**, as stated by SpaceX.



# Conclusions

---



- Other insights obtained from the present analysis were:
- **KSC LC-39A** is the launch site with the highest number of successful missions and **VAFB SLC-4E** is the one with the lowest.
- **ES-L1, GEO, HEO and SSO** are the orbit types with the highest success rates (100%) but with a small number of missions. On the other hand, **VLEO** is the most common orbit type in recent years.
- The **payload range** with the highest launch success rate is from 1952 kg to 5300 kg.
- **FT and B5** are the **booster versions** with the highest launch success rate.

Thank you!

