



IBM Developer  
SKILLS NETWORK

# Winning Space Race with Data Science

Ajay Saxena  
20<sup>th</sup> Aug'23



# Outline

---

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

# Executive Summary

---

## **The following methodologies were used to analyze data:**

- Data Collection using web scraping and SpaceX API
- Exploratory Data Analysis (EDA), including data wrangling, data visualization and interactive visual analytics
- Machine Learning Prediction.

## **Summary of all results**

- It was possible to collect valuable data from public sources
- EDA allowed to identify which features are the best to predict success of launchings
- Machine Learning Prediction showed the best model to predict which characteristics are important to drive this opportunity by the best way, using all collected data.

# Introduction

---

## **Project background and context**

SpaceX is the most successful company of the commercial space age, making space travel affordable. The company 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 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. Based on public information and machine learning models, we are going to predict if SpaceX will reuse the first stage.

## **Questions to be answered**

- How do variables such as payload mass, launch site, number of flights, and orbits affect the success of the first stage landing?
- Does the rate of successful landings increase over the years?
- What is the best algorithm that can be used for binary classification in this case?





Section 1

# Methodology

# Methodology

---

Data collection methodology:

- Using SpaceX Rest API
- Using Web Scrapping from Wikipedia

Performed data wrangling

- Filtering the data
- Dealing with missing values
- Using One Hot Encoding to prepare the data to a binary classification

Performed exploratory data analysis (EDA) using visualization and SQL

Performed interactive visual analytics using Folium and Plotly Dash

Performed predictive analysis using classification models

- Building, tuning and evaluation of classification models to ensure the best results

# Methodology

---

## Executive Summary

- **Data collection methodology**
  - Data from Space X was obtained from 2 sources:
    - Space X API (<https://api.spacexdata.com/v4/rockets/>)
    - WebScraping  
([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))

## Perform data wrangling

Collected data was enriched by creating a landing outcome label based on outcome data after summarizing and analyzing features

Perform exploratory data analysis (EDA) using visualization and SQL

# Data Collection

---

Data collection process involved a combination of API requests from SpaceX REST API and Web Scraping data from a table in SpaceX's Wikipedia entry.

We had to use both of these data collection methods in order to get complete information about the launches for a more detailed analysis.

Data Columns are obtained by using SpaceX REST API:

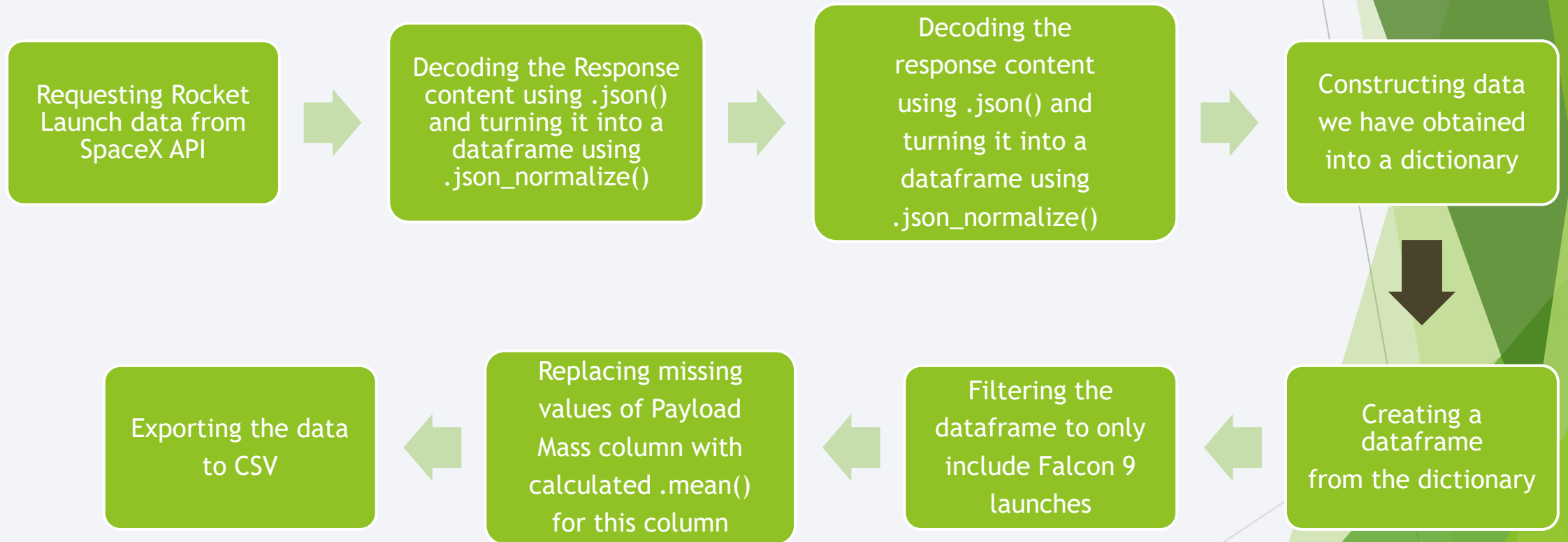
FlightNumber, Date, BoosterVersion, PayloadMass, Orbit, LaunchSite, Outcome, Flights, GridFins, Reused, Legs, LandingPad, Block, ReusedCount, Serial, Longitude, Latitude

Data Columns are obtained by using Wikipedia Web Scraping:

Flight No., Launch site, Payload, PayloadMass, Orbit, Customer, Launch outcome, Version Booster, Booster landing, Date, Time



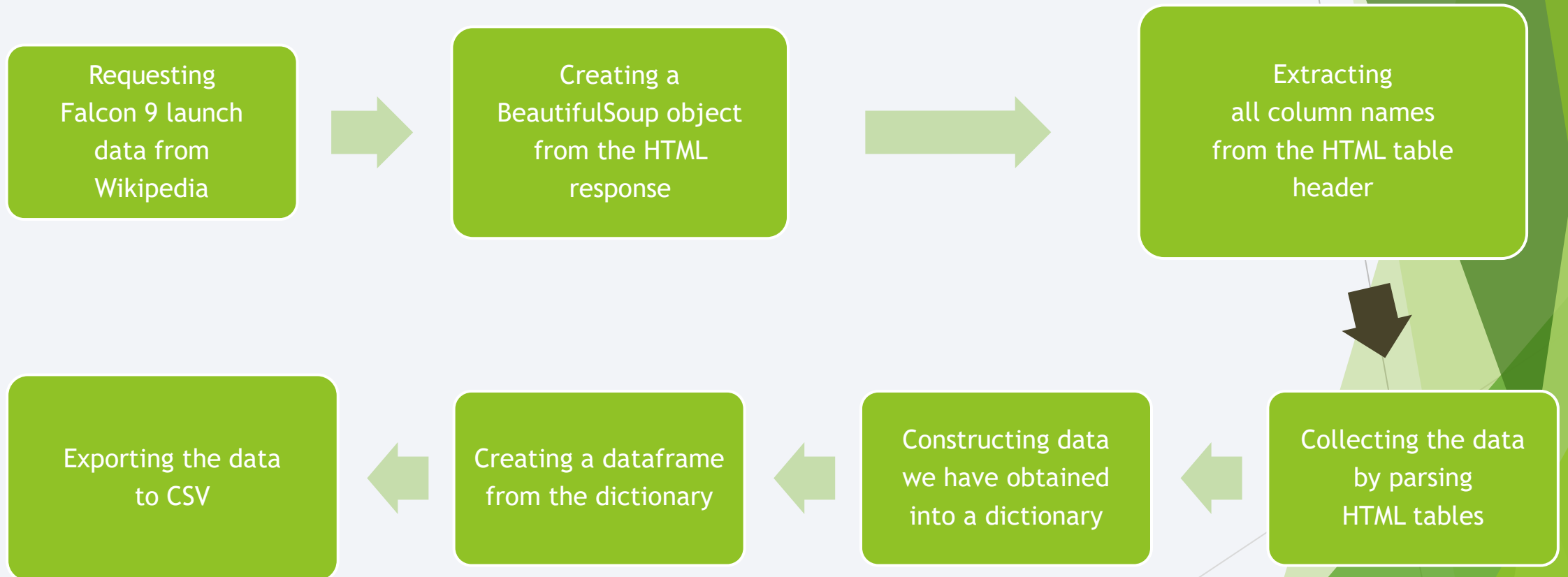
# Data Collection – SpaceX API



[GitHub URL: Data Collection API](#)

# Data Collection - Webscraping

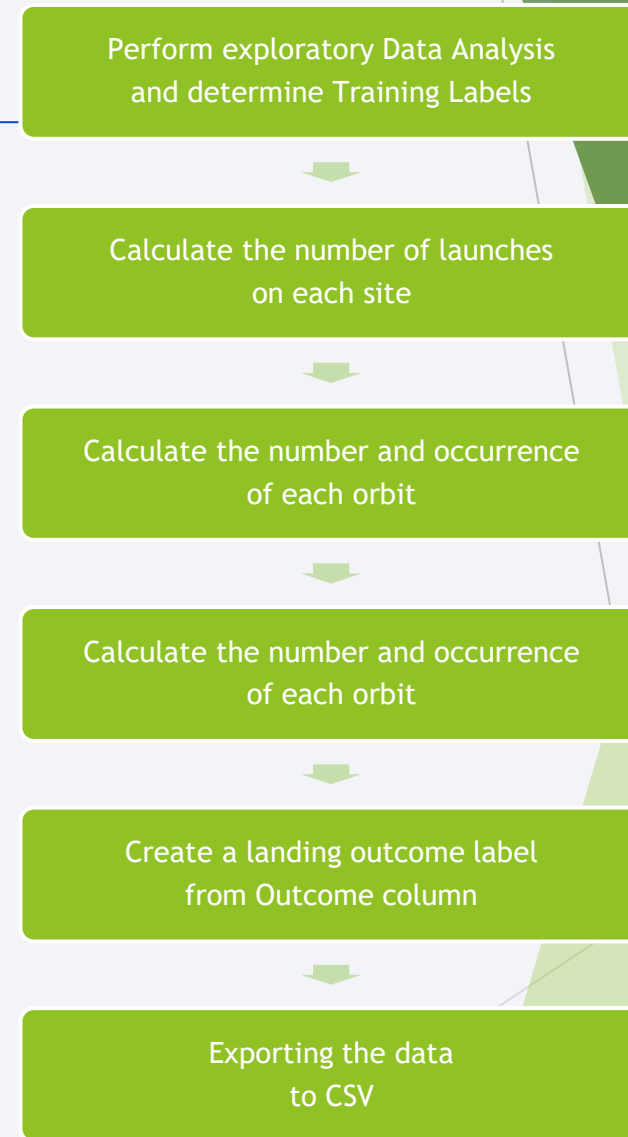
---



[GitHub URL: Data Collection with Webscraping](#)

# Data Wrangling

In the data set, there are several different cases where the booster did not land successfully. Sometimes a landing was attempted but failed due to an accident; for example, True Ocean means the mission outcome was successfully landed to a specific region of the ocean while False Ocean means the mission outcome was unsuccessfully landed to a specific region of the ocean. True RTLS means the mission outcome was successfully landed to a ground pad False RTLS means the mission outcome was unsuccessfully landed to a ground pad. True ASDS means the mission outcome was successfully landed on a drone ship False ASDS means the mission outcome was unsuccessfully landed on a drone ship. We mainly convert those outcomes into Training Labels with “1” means the booster successfully landed, “0” means it was unsuccessful



# EDA with Data Visualization

---

Charts were plotted:

Flight Number vs. Payload Mass, Flight Number vs. Launch Site, Payload Mass vs. Launch Site, Orbit Type vs. Success Rate, Flight Number vs. Orbit Type, Payload Mass vs Orbit Type and Success Rate Yearly Trend

Scatter plots show the relationship between variables. If a relationship exists, they could be used in machine learning model.

Bar charts show comparisons among discrete categories. The goal is to show the relationship between the specific categories being compared and a measured value.

Line charts show trends in data over time (time series).

# EDA with SQL

---

Performed SQL queries:

- Displaying the names of the unique launch sites in the space mission
- Displaying 5 records where launch sites begin with the string 'CCA'
- Displaying the total payload mass carried by boosters launched by NASA (CRS)
- Displaying average payload mass carried by booster version F9 v1.1
- Listing the date when the first successful landing outcome in ground pad was achieved
- Listing the names of the boosters which have success in drone ship and have payload mass greater than 4000 but less than 6000
- Listing the total number of successful and failure mission outcomes
- Listing the names of the booster versions which have carried the maximum payload mass
- Listing the failed landing outcomes in drone ship, their booster versions and launch site names for the months in year 2015
- Ranking 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

[GitHub URL: EDA with SQL](#)



# Build an Interactive Map with Folium

---

Markers of all Launch Sites:

- Added Marker with Circle, Popup Label and Text Label of NASA Johnson Space Center using its latitude and longitude coordinates as a start location.
- Added Markers with Circle, Popup Label and Text Label of all Launch Sites using their latitude and longitude coordinates to show their geographical locations and proximity to Equator and coasts.

Coloured Markers of the launch outcomes for each Launch Site:

- Added coloured Markers of success (Green) and failed (Red) launches using Marker Cluster to identify which launch sites have relatively high success rates.

Distances between a Launch Site to its proximities:

- Added coloured Lines to show distances between the Launch Site KSC LC-39A (as an example) and its proximities like Railway, Highway, Coastline and Closest City

# Build a Dashboard with Plotly Dash

---

Launch Sites Dropdown List:

- Added a dropdown list to enable Launch Site selection.

Pie Chart showing Success Launches (All Sites/Certain Site):

- Added a pie chart to show the total successful launches count for all sites and the Success vs. Failed counts for the site, if a specific Launch Site was selected.

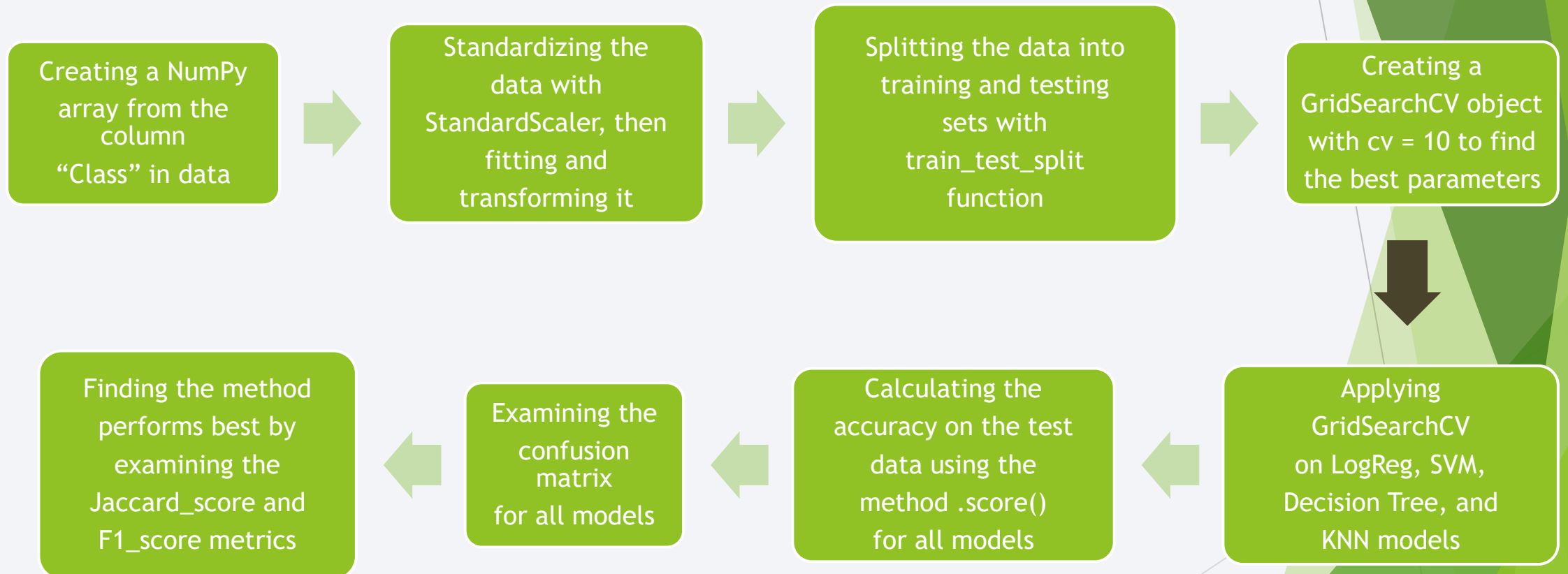
Slider of Payload Mass Range:

- Added a slider to select Payload range.

Scatter Chart of Payload Mass vs. Success Rate for the different Booster Versions:

- Added a scatter chart to show the correlation between Payload and Launch Success

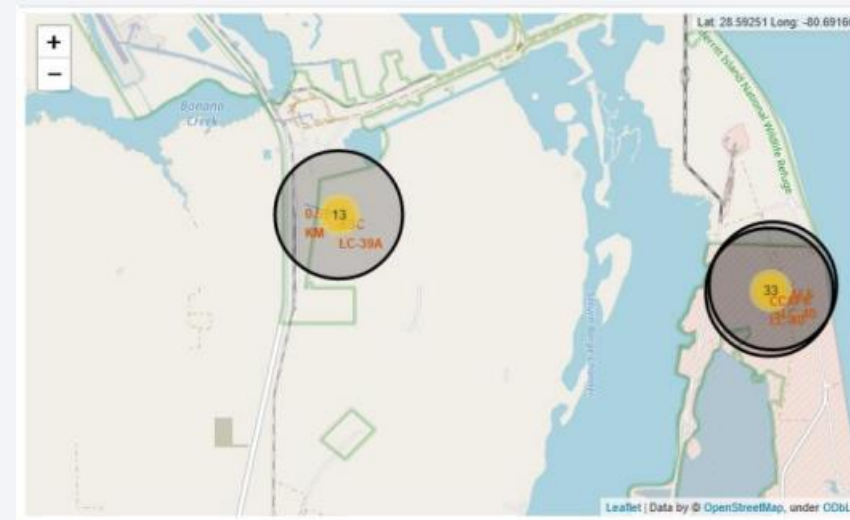
# Predictive Analysis (Classification)



[GitHub URL: Machine Learning Prediction](#)

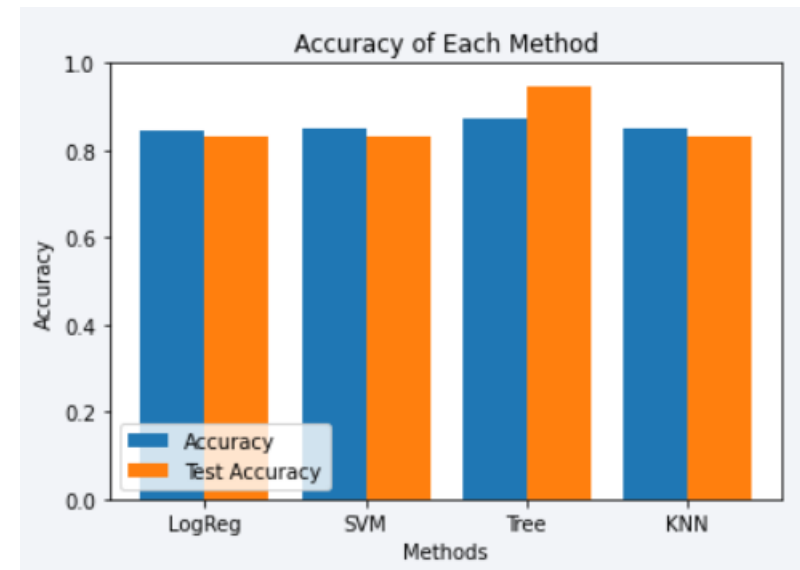
# Results

- Using interactive analytics was possible to identify that launch sites use to be in safety places, near sea, for example and have a good logistic infrastructure around.
- Most launches happens at east cost launch sites.



# Results

- Predictive Analysis showed that Decision Tree Classifier is the best model to predict successful landings, having accuracy over 87% and accuracy for test data over 94%.



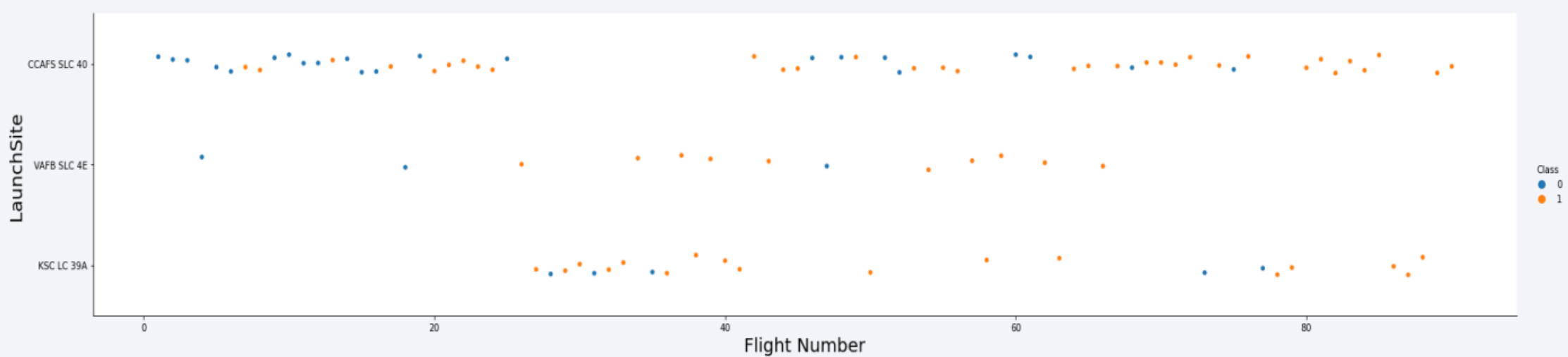




Section 2

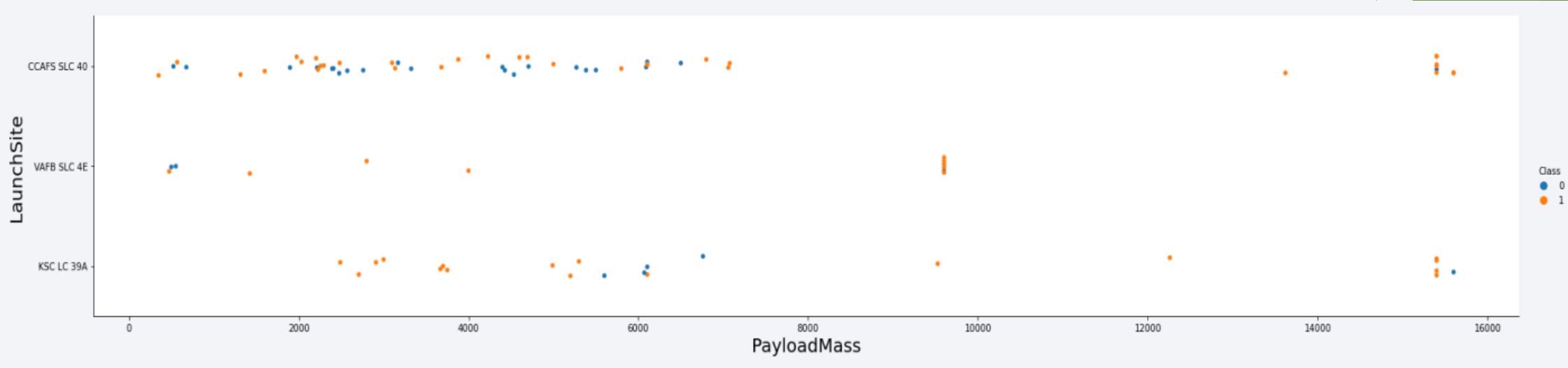
# Insights drawn from EDA

# Flight Number vs. Launch Site



- According to the plot above, it's possible to verify that the best launch site nowadays is CCAFS SLC 40, where most of recent launches were successful;
- In second place VAFB SLC 4E and third place KSC LC 39A;
- It's also possible to see that the general success rate improved over time.

# Payload vs. Launch Site



Payloads over 9,000kg (about the weight of a school bus) have excellent success rate

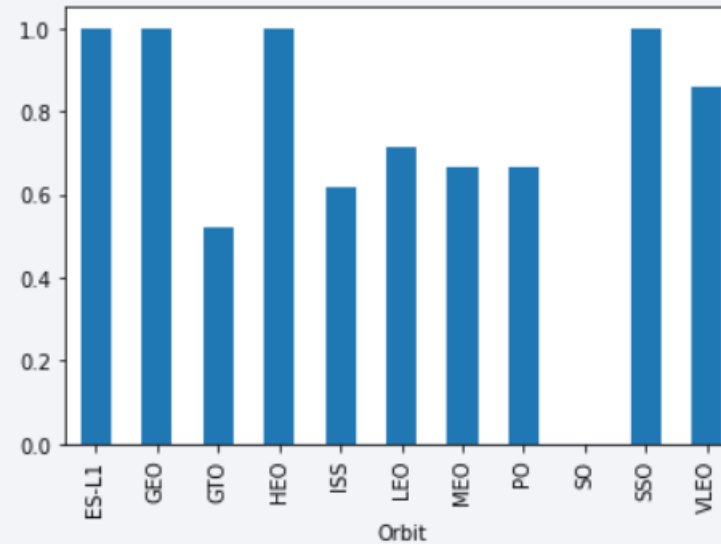
- Payloads over 12,000kg seems to be possible only on CCAFS SLC 40 and KSC LC 39A launch sites.



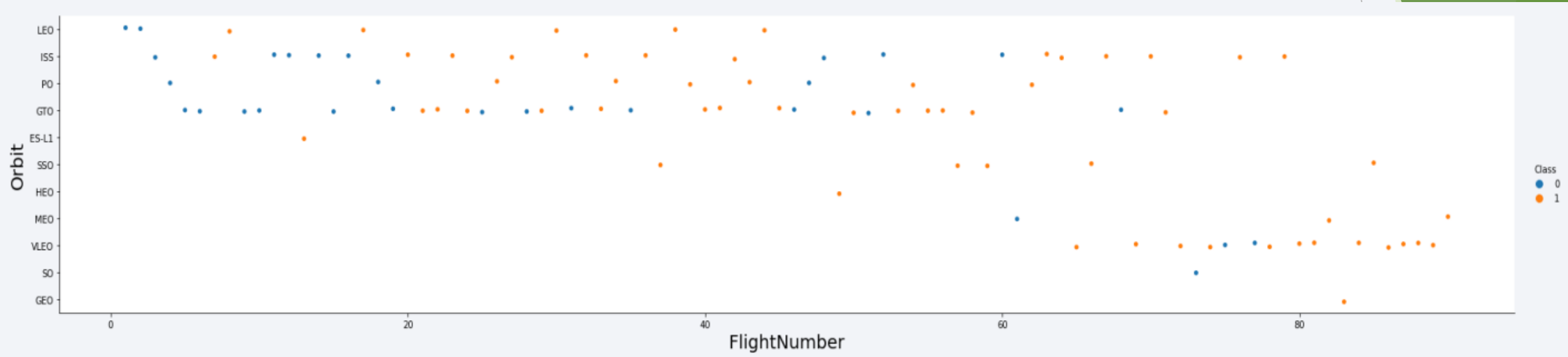
# Success Rate vs. Orbit Type

The biggest success rates happens to orbits:

- ES-L1;
- GEO;
- HEO; and
- SSO.
- Followed by:
- VLEO (above 80%); and
- LFO (above 70%).



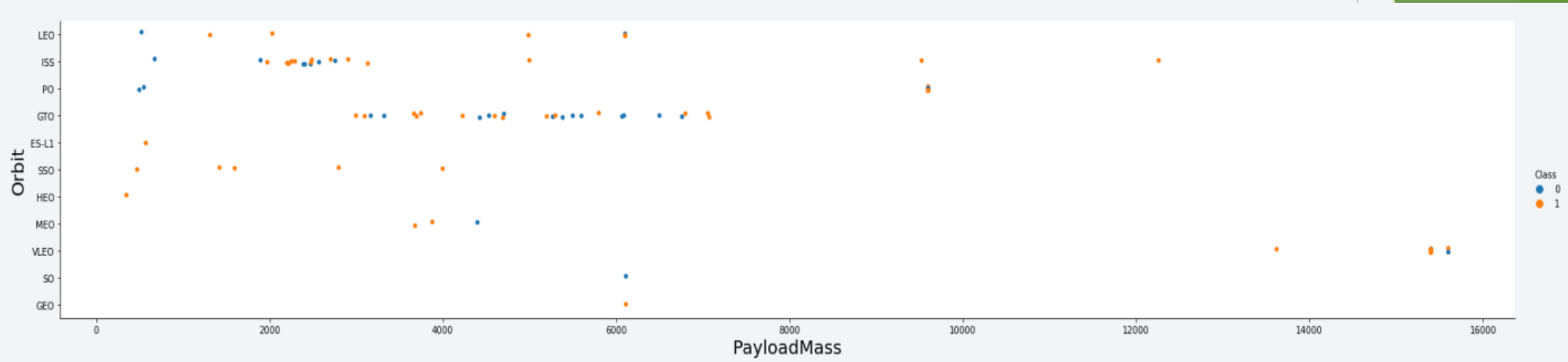
# Flight Number vs. Orbit Type



- ▶ Apparently, success rate improved over time to all orbits
- ▶ VLEO orbit seems a new business opportunity, due to recent increase of its frequency.



# Payload vs. Orbit Type

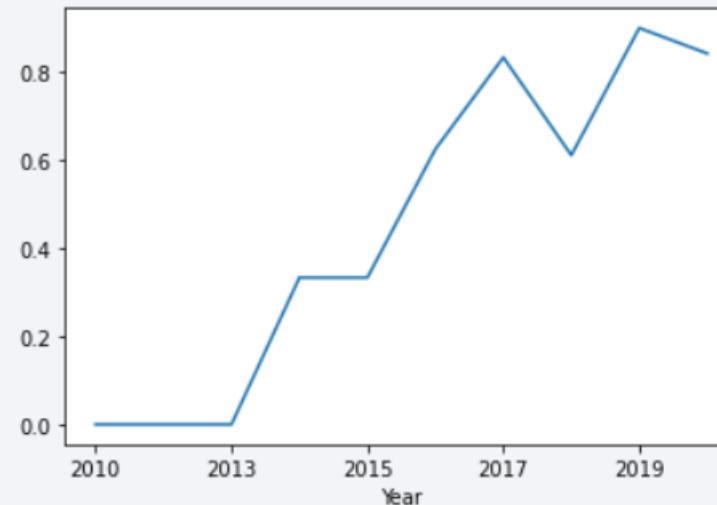


Apparently, there is no relation between payload and success rate to orbit GTO;  
ISS orbit has the widest range of payload and a good rate of success;  
There are few launches to the orbits SO and GEO.

# Launch Success Yearly Trend

---

- ▶ Success rate started increasing in 2013 and kept until 2020;
- ▶ It seems that the first three years were a period of adjusts and improvement of technology.



# All Launch Site Names

---

- ▶ According to data, there are four launch sites:

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

- ▶ They are obtained by selecting unique occurrences of “launch\_site” values from the dataset.

# Launch Site Names Begin with 'CCA'

- They are obtained by selecting unique occurrences of “launch\_site” values from the dataset.

Date	Time UTC	Booster Version	Launch Site	Payload	Payload Mass kg	Orbit	Customer	Mission Outcome	Landing Outcome
2010-06-04	18:45:00	F9 v1.0 B0003	CCAFS LC-40	Dragon Spacecraft Qualification Unit	0	LEO	SpaceX	Success	Failure (parachute)
2010-12-08	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
2012-10-08	00:35:00	F9 v1.0 B0006	CCAFS LC-40	SpaceX CRS-1	500	LEO (ISS)	NASA (CRS)	Success	No attempt
2013-03-01	15:10:00	F9 v1.0 B0007	CCAFS LC-40	SpaceX CRS-2	677	LEO (ISS)	NASA (CRS)	Success	No attempt

- Here we can see five samples of Cape Canaveral launches.

# Total Payload Mass

---

- ▶ Total payload carried by boosters from NASA:

Total Payload (kg)
111.268

- ▶ Total payload calculated above, by summing all payloads whose codes contain 'CRS', which corresponds to NASA.



# Average Payload Mass by F9 v1.1

---

- ▶ Average payload mass carried by booster version F9 v1.1:

Avg Payload (kg)
2.928

- ▶ Filtering data by the booster version above and calculating the average payload mass we obtained the value of 2,928 kg.

# First Successful Ground Landing Date

---

- ▶ First successful landing outcome on ground pad:

Min Date
2015-12-22

- ▶ By filtering data by successful landing outcome on ground pad and getting the minimum value for date it's possible to identify the first occurrence, that happened on 12/22/2015.

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

---

- ▶ Boosters which have successfully landed on drone ship and had payload mass greater than 4000 but less than 6000

Booster Version
F9 FT B1021.2
F9 FT B1031.2
F9 FT B1022
F9 FT B1026

- ▶ Selecting distinct booster versions according to the filters above, these 4 are the result.

# Total Number of Successful and Failure Mission Outcomes

---

- ▶ Number of successful and failure mission outcomes:

Mission Outcome	Occurrences
Success	99
Success (payload status unclear)	1
Failure (in flight)	1

- ▶ Grouping mission outcomes and counting records for each group led us to the summary above.

# Boosters Carried Maximum Payload

- ▶ Boosters which have carried the maximum payload mass

Booster Version (...)
F9 B5 B1048.4
F9 B5 B1048.5
F9 B5 B1049.4
F9 B5 B1049.5
F9 B5 B1049.7
F9 B5 B1051.3

Booster Version
F9 B5 B1051.4
F9 B5 B1051.6
F9 B5 B1056.4
F9 B5 B1058.3
F9 B5 B1060.2
F9 B5 B1060.3

- ▶ These are the boosters which have carried the maximum payload mass registered in the dataset.

# 2015 Launch Records

---

- ▶ Failed landing outcomes in drone ship, their booster versions, and launch site names for in year 2015

Booster Version	Launch Site
F9 v1.1 B1012	CCAFS LC-40
F9 v1.1 B1015	CCAFS LC-40

- ▶ The list above has the only two occurrences.



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

---

- ▶ Ranking of all landing outcomes between the date 2010-06-04 and 2017-03-20:

Landing Outcome	Occurrences
No attempt	10
Failure (drone ship)	5
Success (drone ship)	5
Controlled (ocean)	3
Success (ground pad)	3
Failure (parachute)	2
Uncontrolled (ocean)	2
Precluded (drone ship)	1

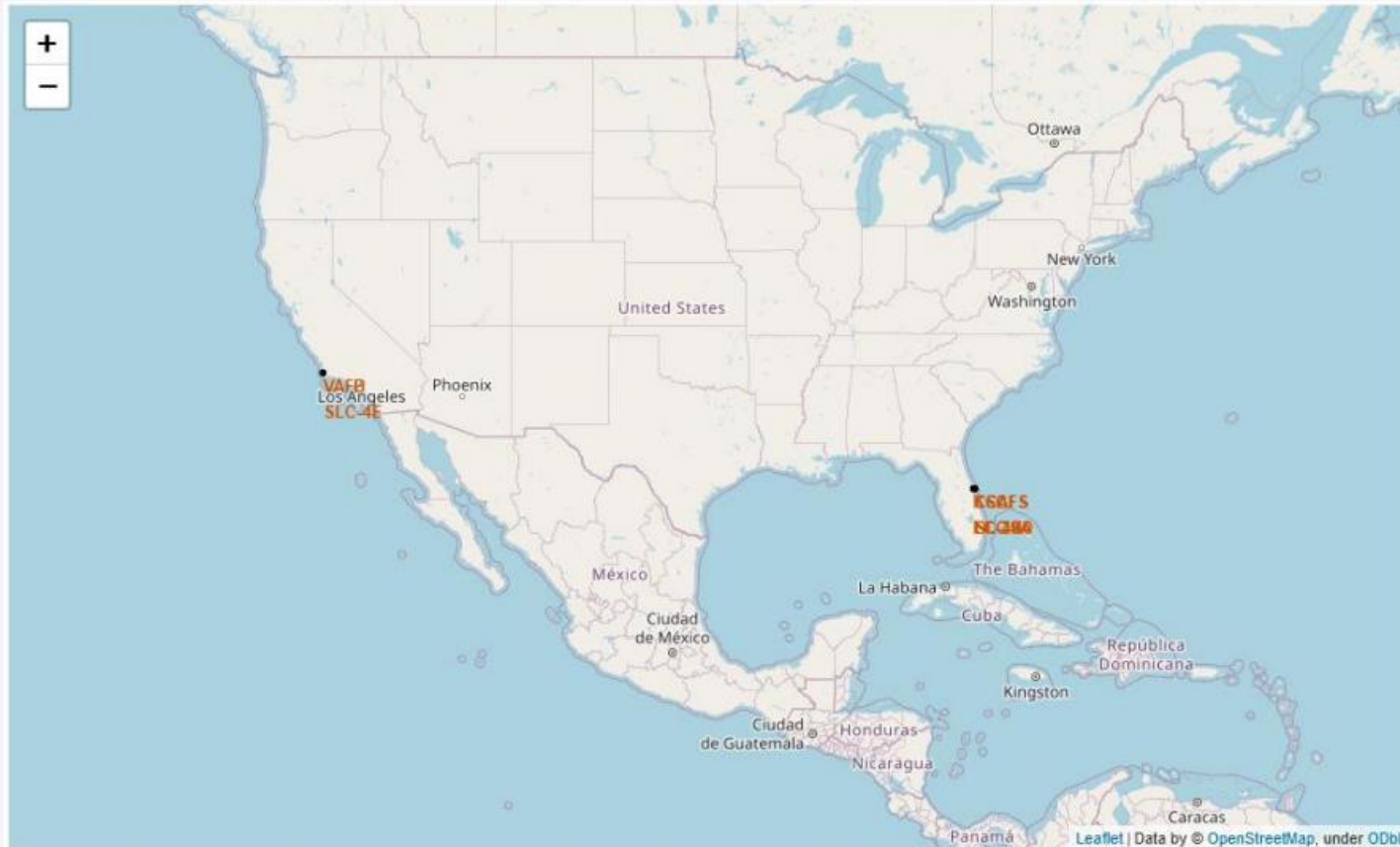
- ▶ This view of data alerts us that “No attempt” must be taken in account.

A satellite view of Earth from space, showing the curvature of the planet and city lights at night. The background is a deep blue space with stars. The Earth's surface is a mix of dark blue oceans and bright yellow city lights. The text is overlaid on the left side of the image.

Section 3

# Launch Sites Proximities Analysis

# All launch sites



- Launch sites are near sea, probably by safety, but not too far from roads and railroads.

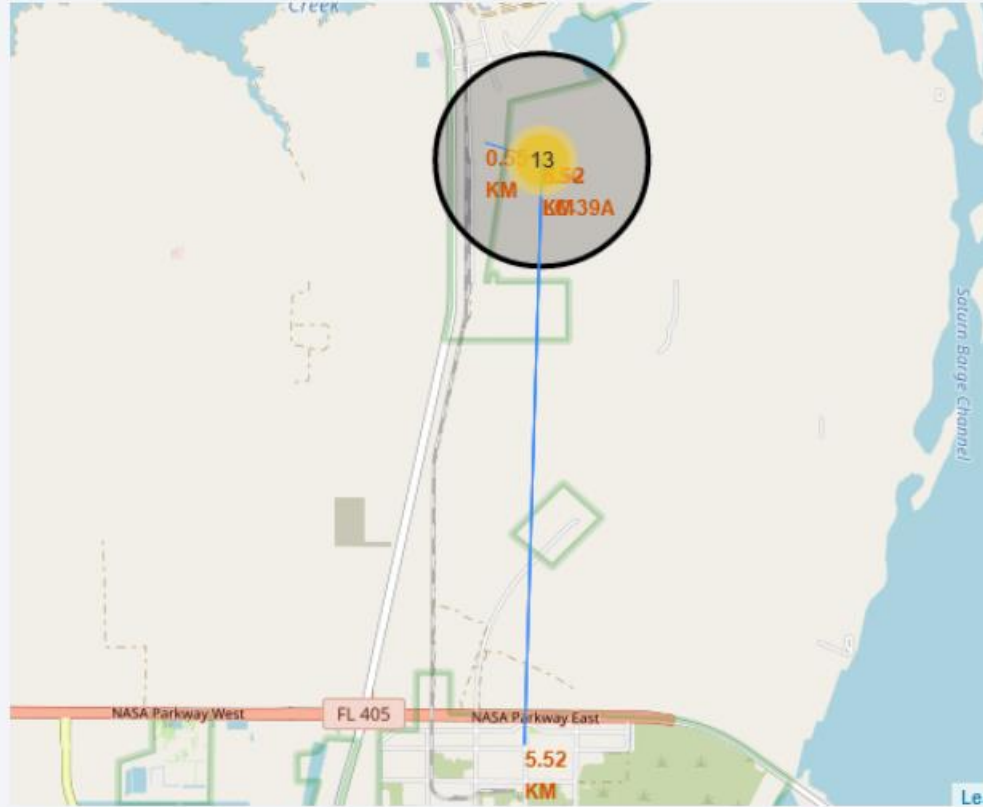
# Launch Outcomes by Site

- ▶ Example of KSC LC-39A launch site launch outcomes



- ▶ Green markers indicate successful and red ones indicate failure.

# Logistics and Safety



- Launch site KSC LC-39A has good logistics aspects, being near railroad and road and relatively far from inhabited areas.



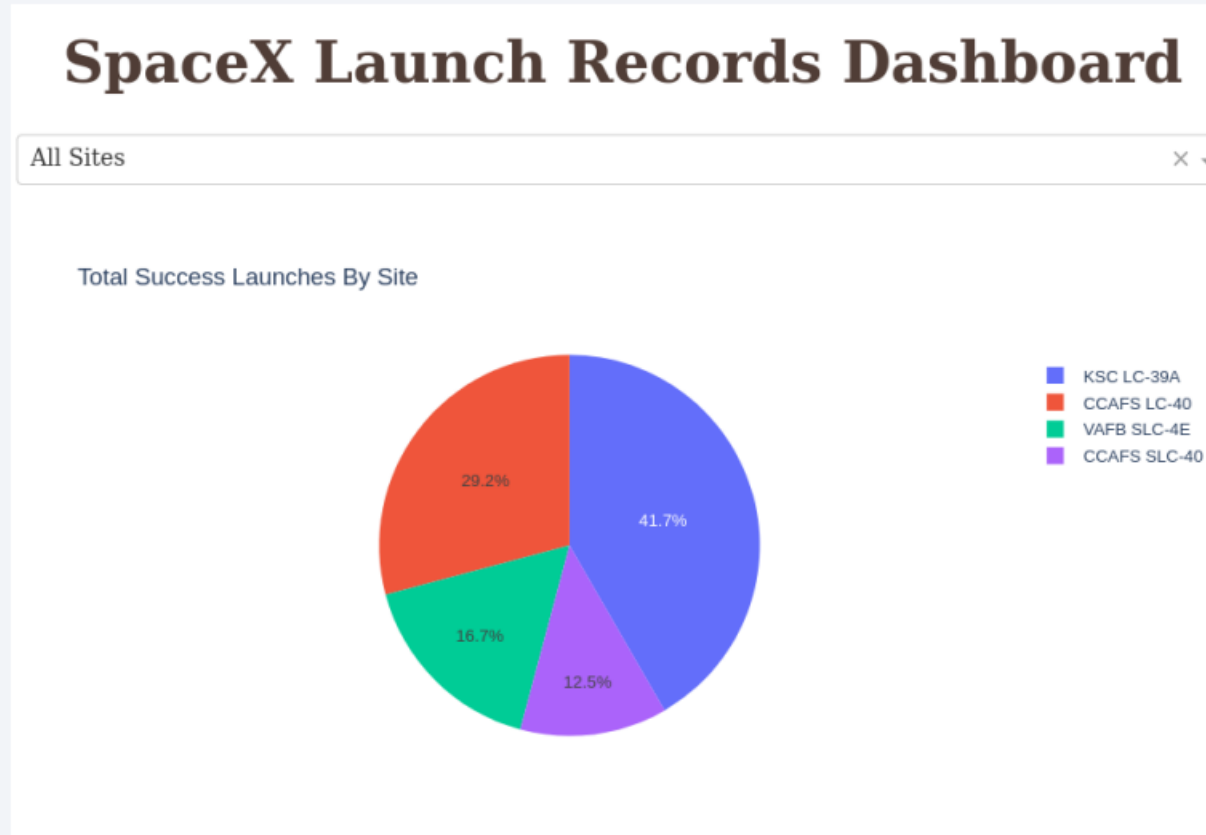


Section 4

# Build a Dashboard with Plotly Dash



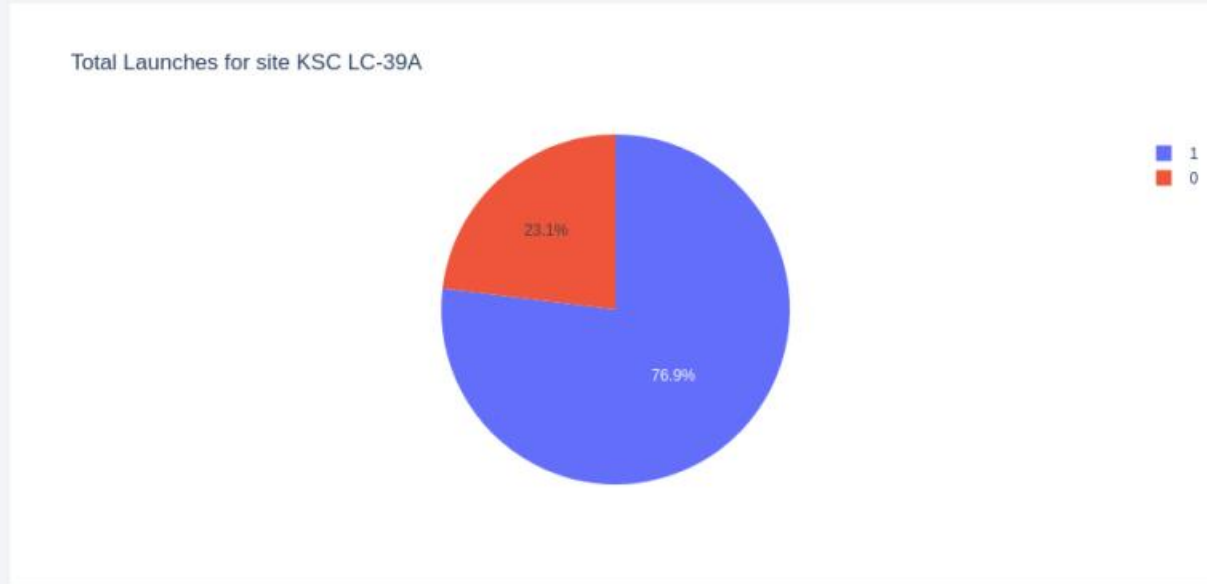
# Successful Launches by Site



- The place from where launches are done seems to be a very important factor of success of missions.

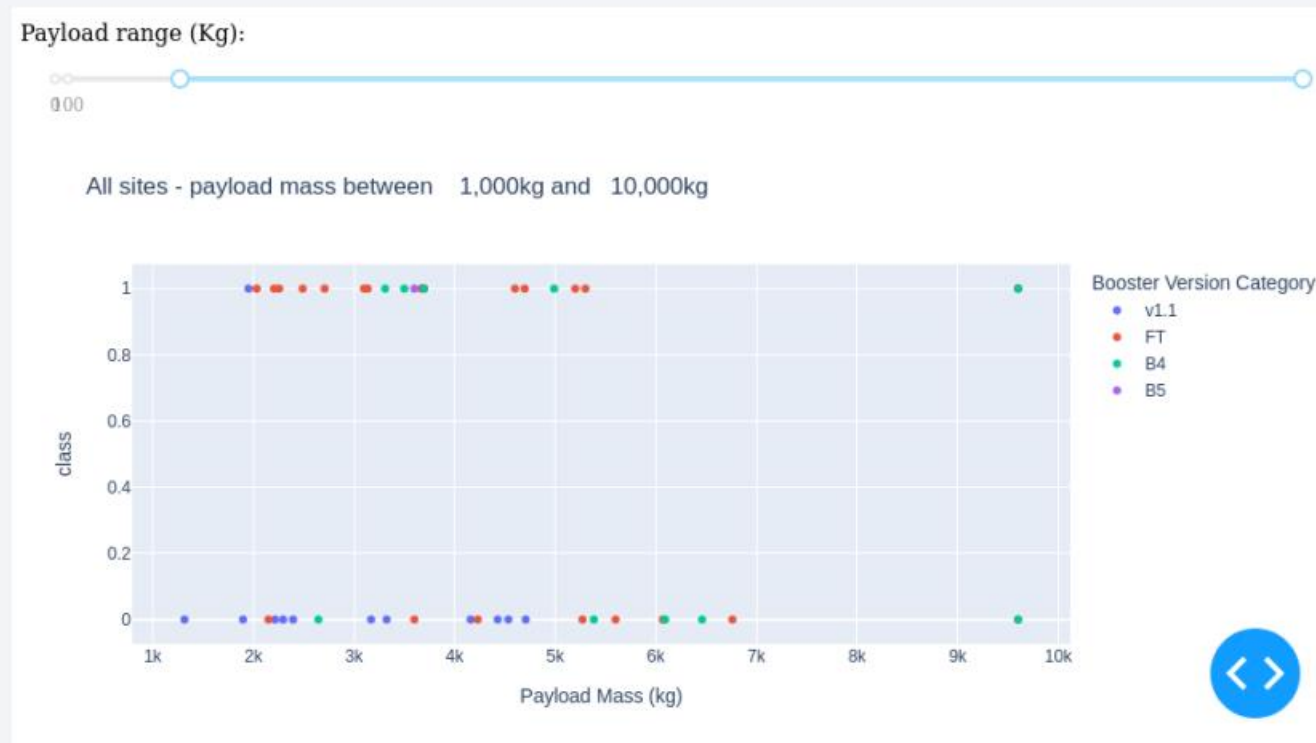
# Launch Success Ratio for KSC LC-39A

---



- Launch Success Ratio for KSC LC-39A

# Payload vs. Launch Outcome



- Payloads under 6,000kg and FT boosters are the most successful combination.

# Payload vs. Launch Outcome



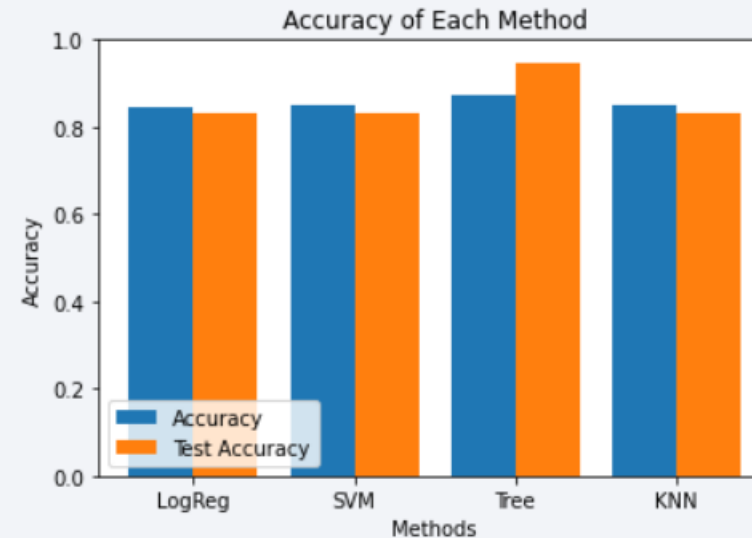
► There's not enough data to estimate risk of launches over 7,000kg

Section 5

# Predictive Analysis (Classification)

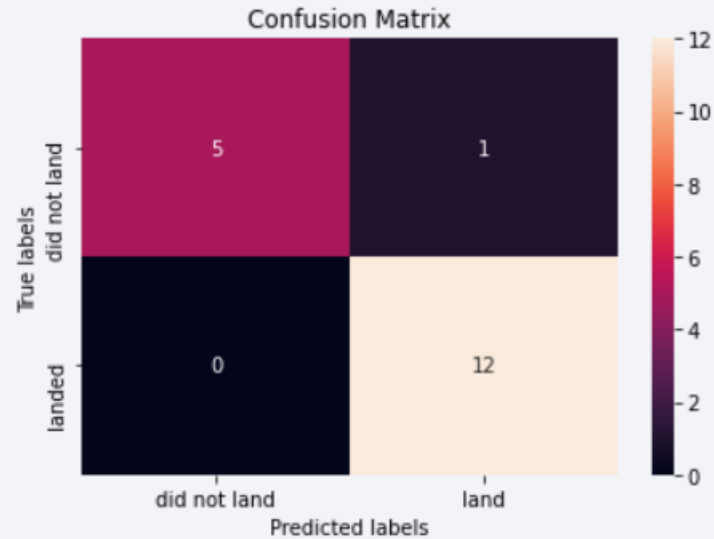
# Classification Accuracy

- ▶ Four classification models were tested, and their accuracies are plotted beside
- ▶ The model with the highest classification accuracy is Decision Tree Classifier, which has accuracies over than 87%.





# Confusion Matrix of Decision Tree Classifier



- Confusion matrix of Decision Tree Classifier proves its accuracy by showing the big numbers of true positive and true negative compared to the false ones.

# Conclusions

---

- ▶ Different data sources were analyzed, refining conclusions along the process
- ▶ The best launch site is KSC LC-39A
- ▶ Launches above 7,000kg are less risky
- ▶ Although most of mission outcomes are successful, successful landing outcomes seem to improve over time, according the evolution of processes and rockets
- ▶ Decision Tree Classifier can be used to predict successful landings and increase profits.

# Appendix

---

- ▶ As an improvement for model tests, it's important to set a value to `np.random.seed` variable
- ▶ Folium didn't show maps on Github, so I took screenshots.

Thank you!

