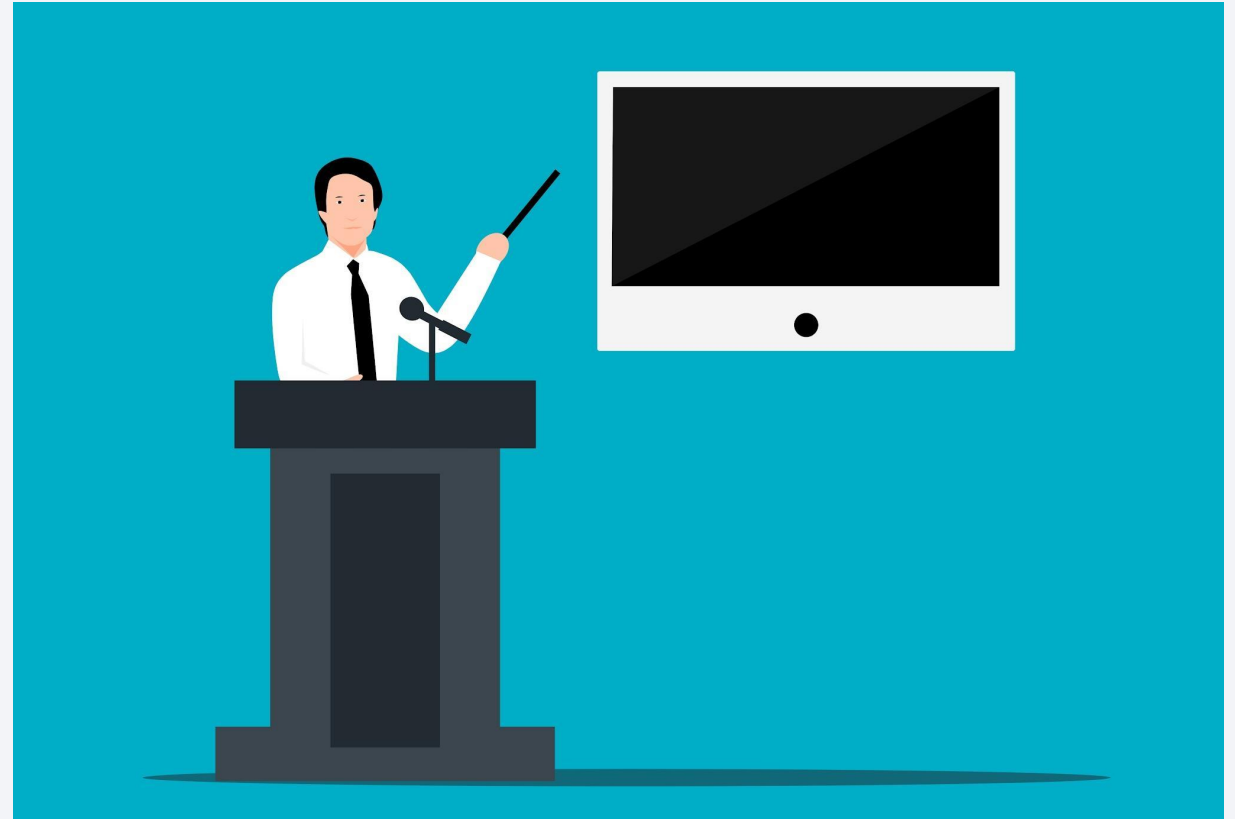# Winning Space Race with Data Science

Parshv Patel
6 January 2025

# Outline

- [Executive Summary](#)

- [Introduction](#)

- [Methodology](#)

- [Results](#)

- [Conclusion](#)

- [Appendix](#)

# Executive Summary

**Summary of methodologies:**

- Data collection using Rest API and Web Scraping
- Data wrangling and preparation
- Exploratory Data Analysis with Data Visualization
- Exploratory Data Analysis with SQL
- Building an interactive map with Folium
- Building a Dashboard with Plotly Dash
- Modeling Classification Predictive analysis

**Summary of all results:**

- Exploratory Data Analysis results
- Interactive analytics demo in screenshots
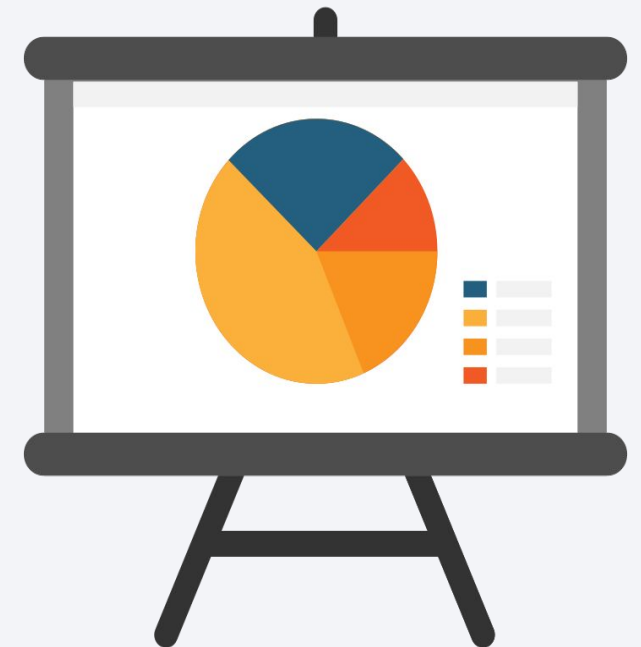- Model evaluation and Predictive analysis results

3

# 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 Scraping from Wikipedia

**Perform data wrangling:**

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

**Perform exploratory data analysis (EDA) using visualization and SQL**

**Perform interactive visual analytics using Folium and Plotly Dash**

**Perform predictive analysis using classification models:**

- Modeling different types of classification, hyperparameter tuning using GridSerach, and evaluating the model using accuracy, F1-score, and jaccard index to ensure the best results.

# 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.
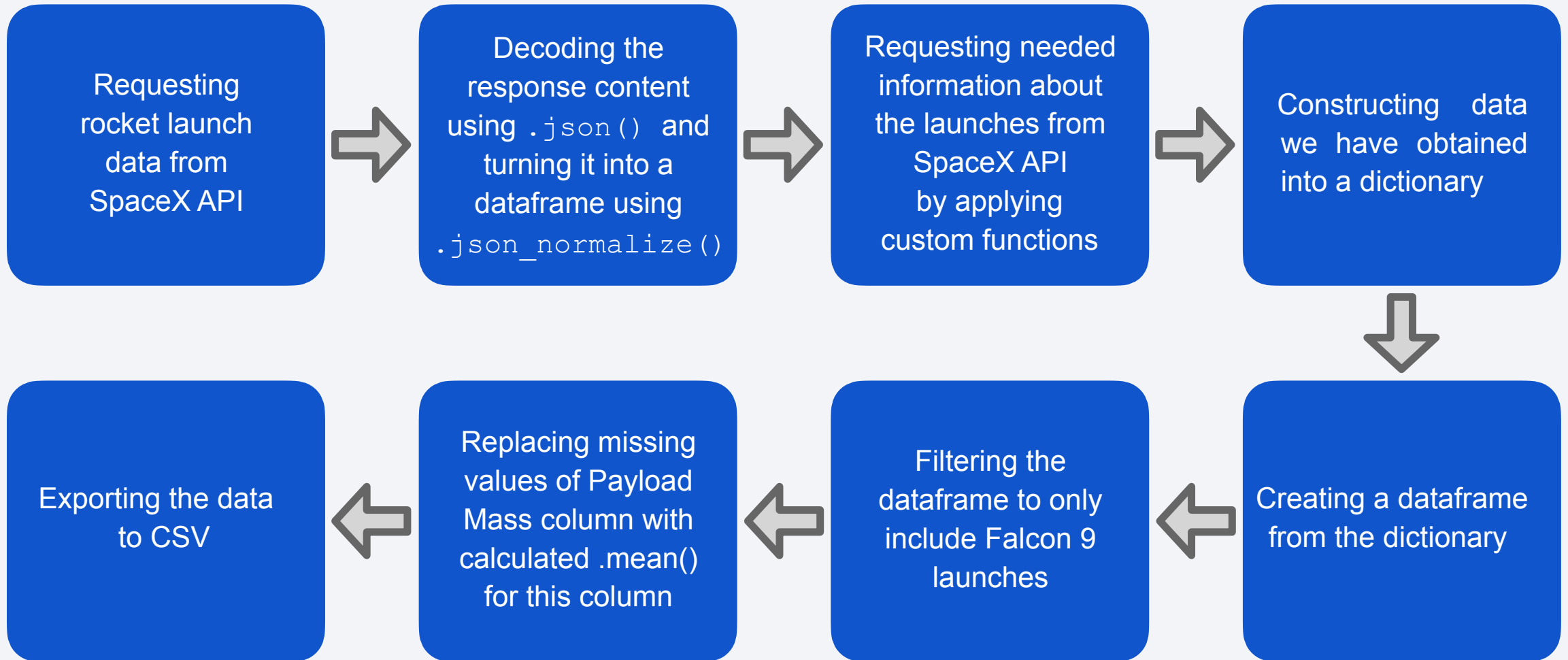
**Data columns are obtained by using SpaceX REST API:**

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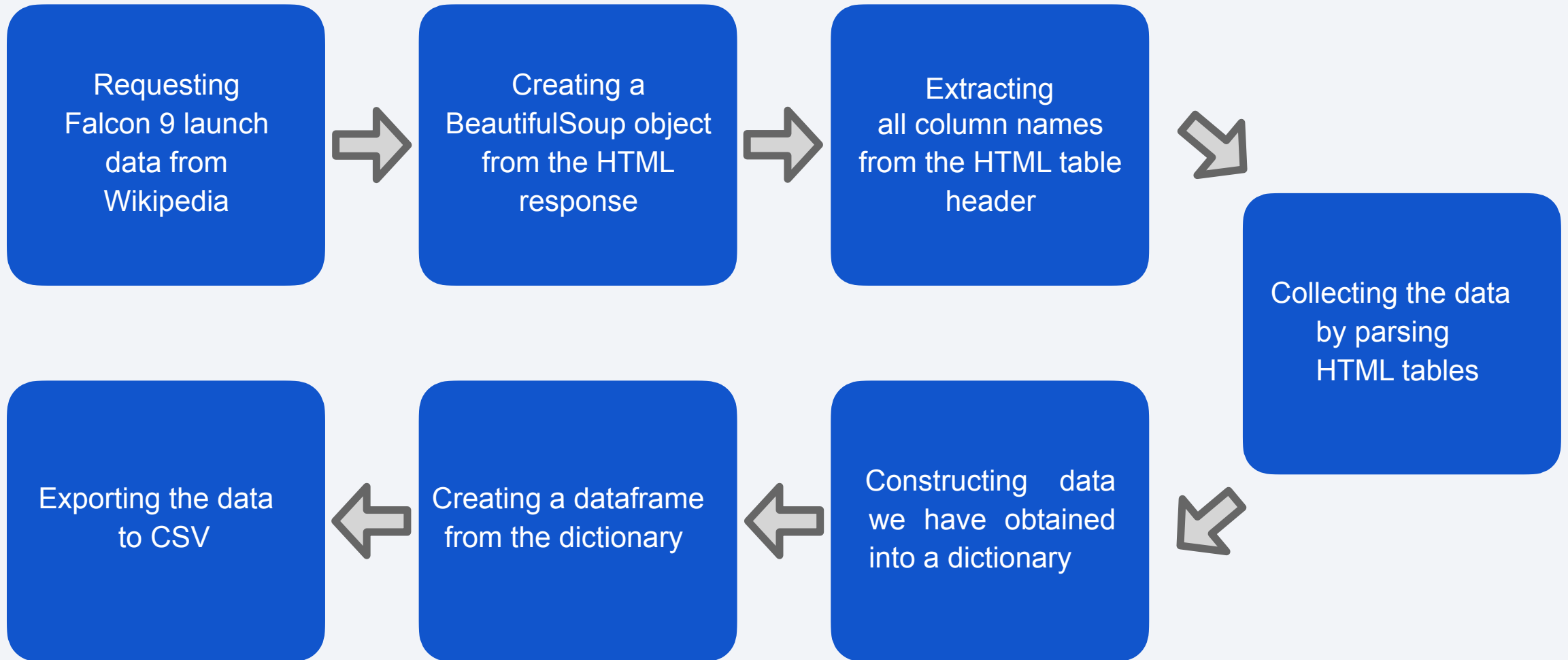
    Flight No., Launch site, Payload, PayloadMass, Orbit, Customer, Launch outcome, Version Booster, Booster landing, Date, Time

# Data Collection- SpaceX API

Requesting rocket launch data from SpaceX API

Decoding the response content using `.json()` and turning it into a dataframe using `.json_normalize()`

Requesting needed information about the launches from SpaceX API by applying custom functions

Constructing data we have obtained into a dictionary

Creating a dataframe from the dictionary

Filtering the dataframe to only include Falcon 9 launches

Replacing missing values of Payload Mass column with calculated .mean() for this column

Exporting the data to CSV

GitHub URL: Data Collection API

# Data Collection - Web Scraping

Requesting
Falcon 9 launch
data from
Wikipedia

Creating a
BeautifulSoup object
from the HTML
response

Extracting
all column names
from the HTML table
header

Collecting the data
by parsing
HTML tables

Exporting the data
to CSV

Creating a dataframe
from the dictionary
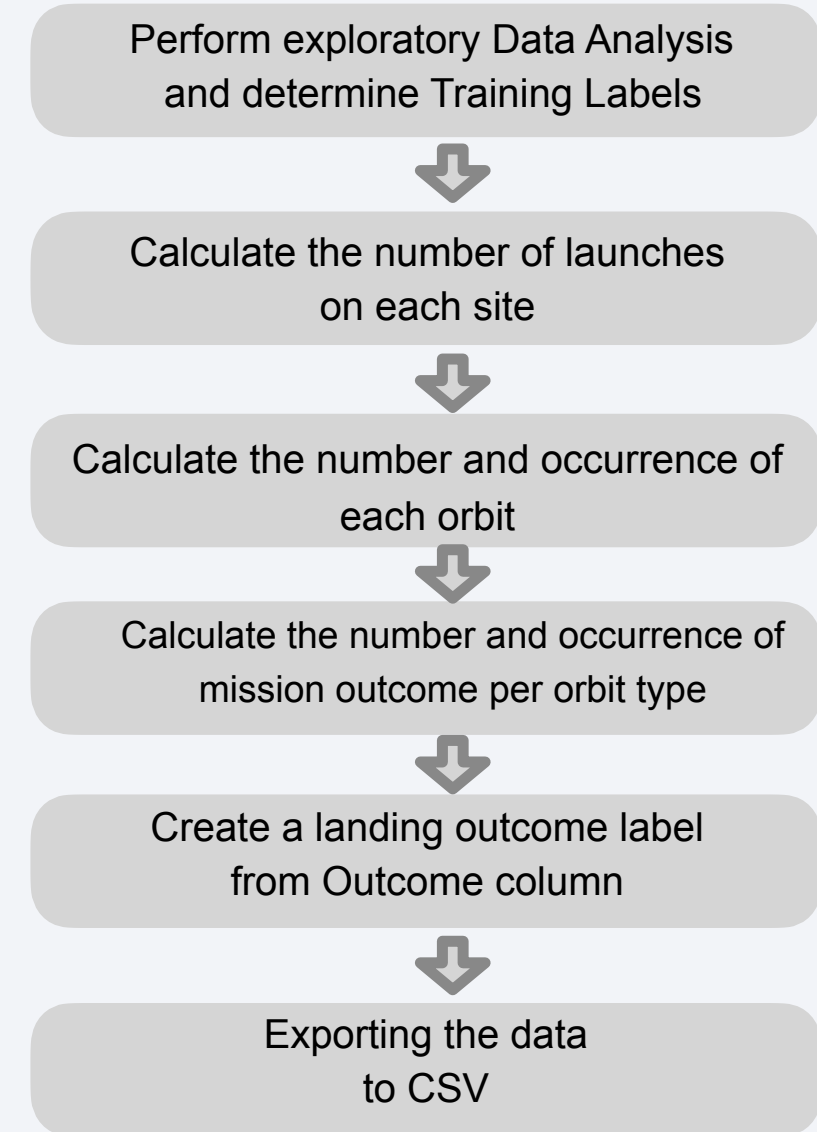
Constructing    data
we  have  obtained
into a dictionary

GitHub URL: Data Collection with Web Scraping

# 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.

GitHub URL: Data Wrangling

Perform exploratory Data Analysis and determine Training Labels

⬇

Calculate the number of launches on each site

⬇

Calculate the number and occurrence of each orbit

⬇

Calculate the number and occurrence of mission outcome per orbit type

⬇

Create a landing outcome label from Outcome column

⬇

Exporting the data to CSV

# EDA with Data Visualization

Charts 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).

GitHub URL: EDA with Data Visualization

# 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.

GitHub URL: Interactive Visual Analytics with Folium

# 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.
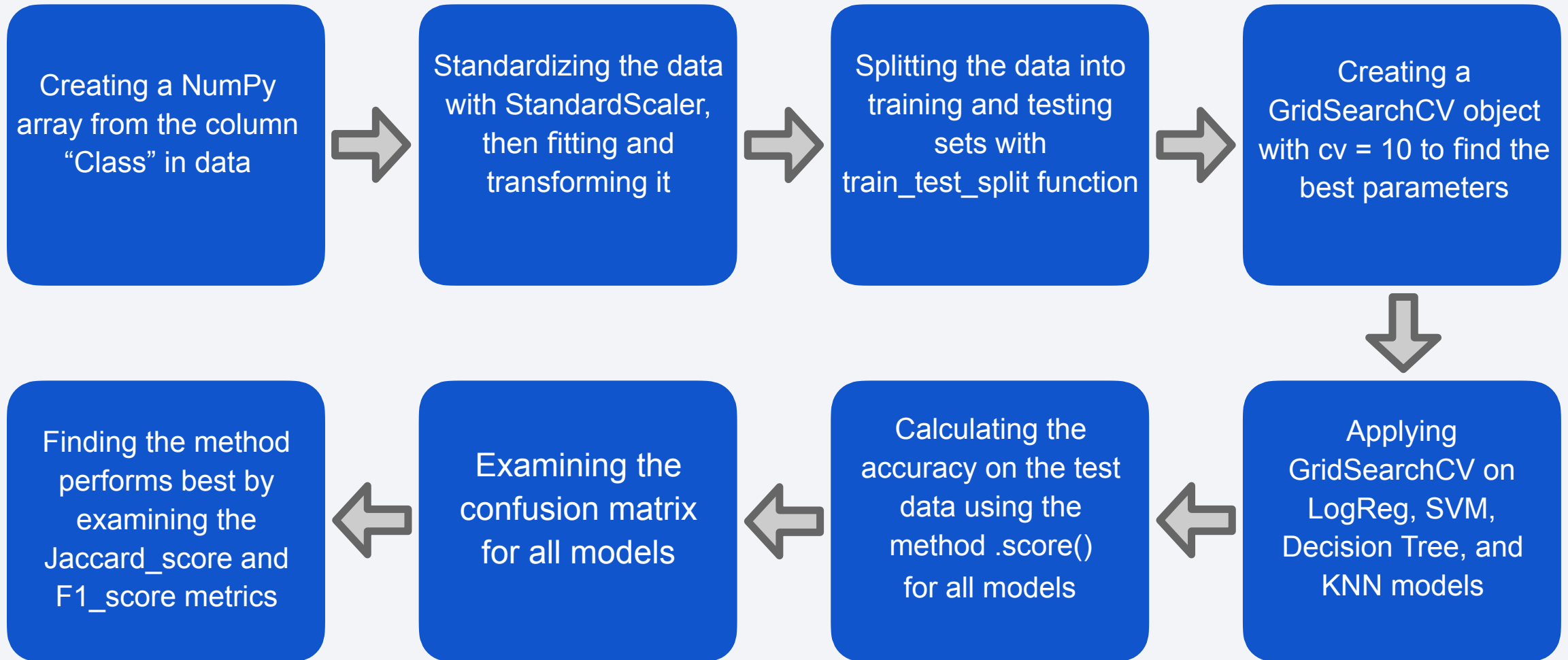
Slider of Payload Mass Range:

- 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.

GitHub URL: SpaceX Dash Application

# Predictive Analysis (Classification)

Creating a NumPy array from the column "Class" in data

→

Standardizing the data with StandardScaler, then fitting and transforming it

→

Splitting the data into training and testing sets with train_test_split function

→

Creating a GridSearchCV object with cv = 10 to find the best parameters

↓

Finding the method performs best by examining the Jaccard_score and F1_score metrics

←

Examining the confusion matrix for all models

←

Calculating the accuracy on the test data using the method .score() for all models

←

Applying GridSearchCV on LogReg, SVM, Decision Tree, and KNN models

GitHub URL: Machine Learning Prediction

# Results

- Exploratory data analysis results

- Interactive analytics demo in screenshots

- Predictive analysis results



publicdomainvectors.org

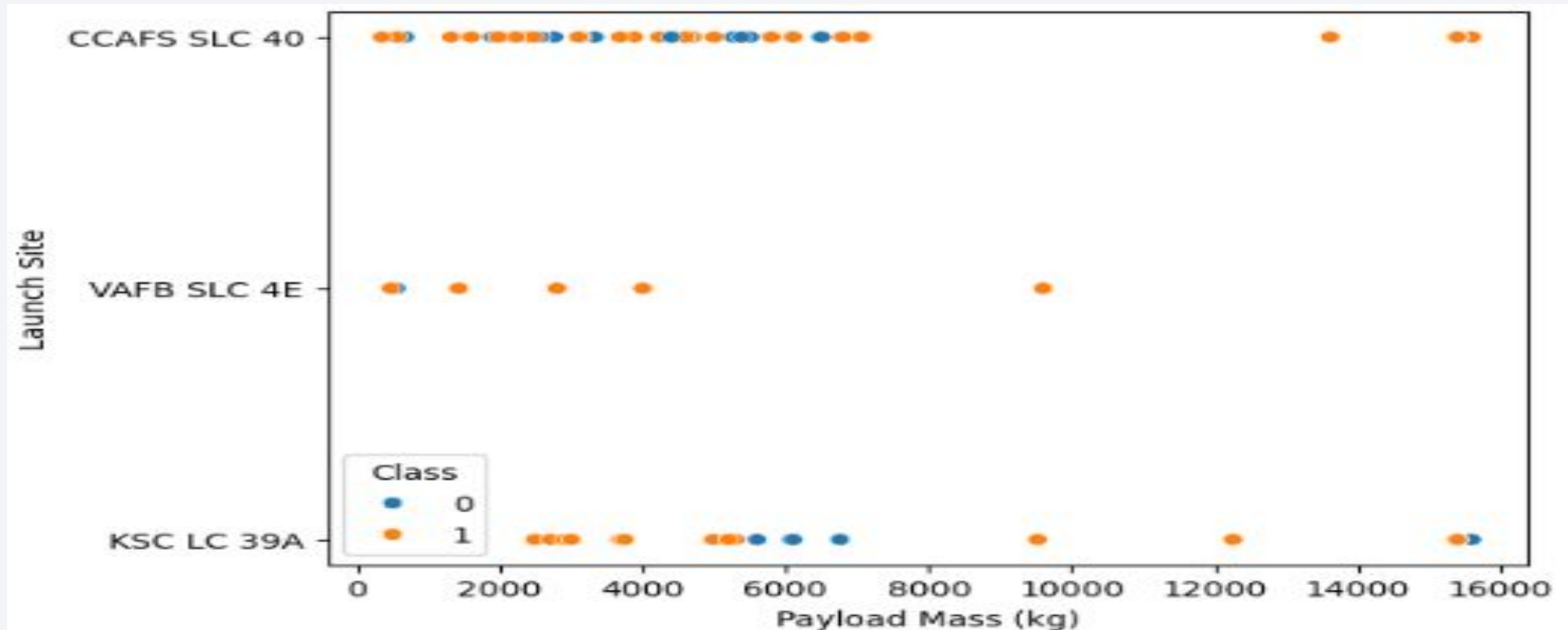Section 2

# Insights drawn from EDA

# EDA with Visualization

# Flight Number vs. Launch Site



- The earliest flights all failed while the latest flights all succeeded.
- The CCAFS SLC 40 launch site has about a half of all launches and higher chance of failure when the flight number was low.
- VAFB SLC 4E and KSC LC 39A have higher success rates.
- It can be assumed that each new launch has a higher rate of success.
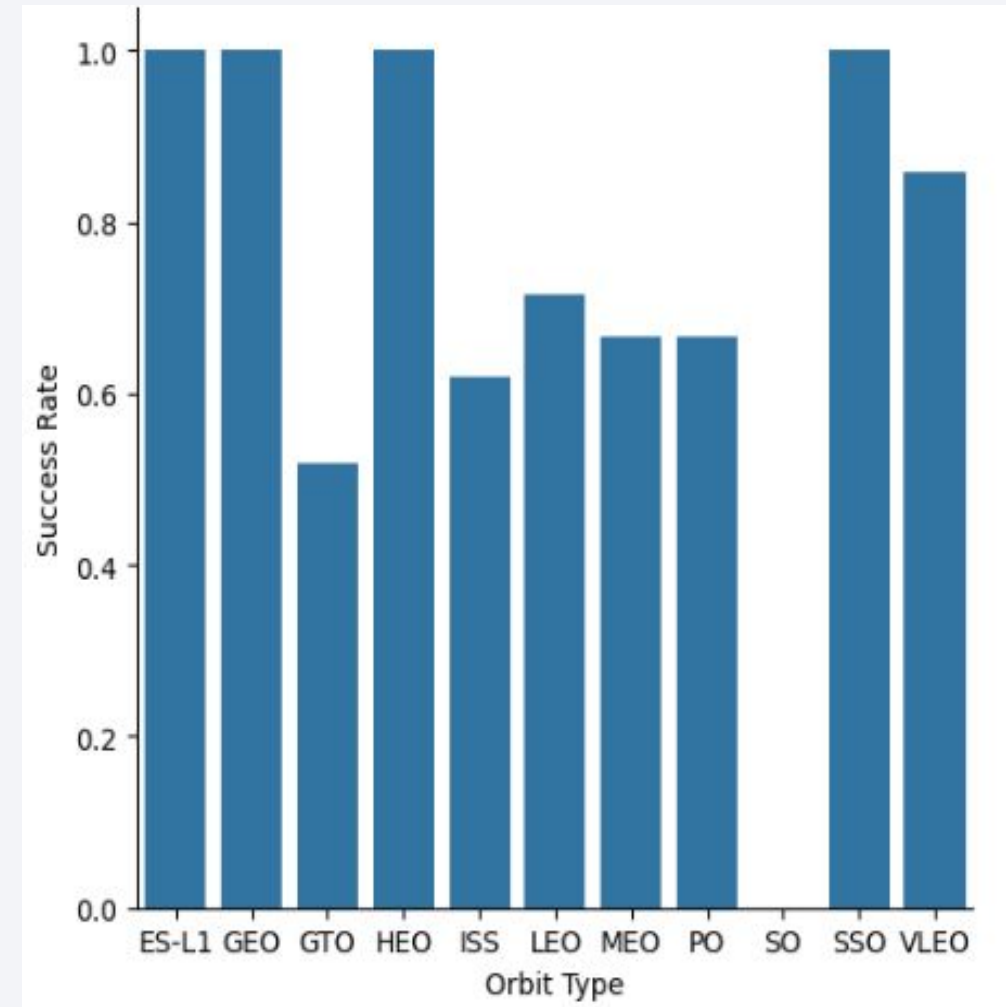
# Payload Mass vs. Launch Site



- For every launch site the higher the payload mass, the higher the success rate.
- Most of the launches with payload mass over 7000 kg were successful.
- KSC LC 39A has a 100% success rate for payload mass under 5500 kg too.
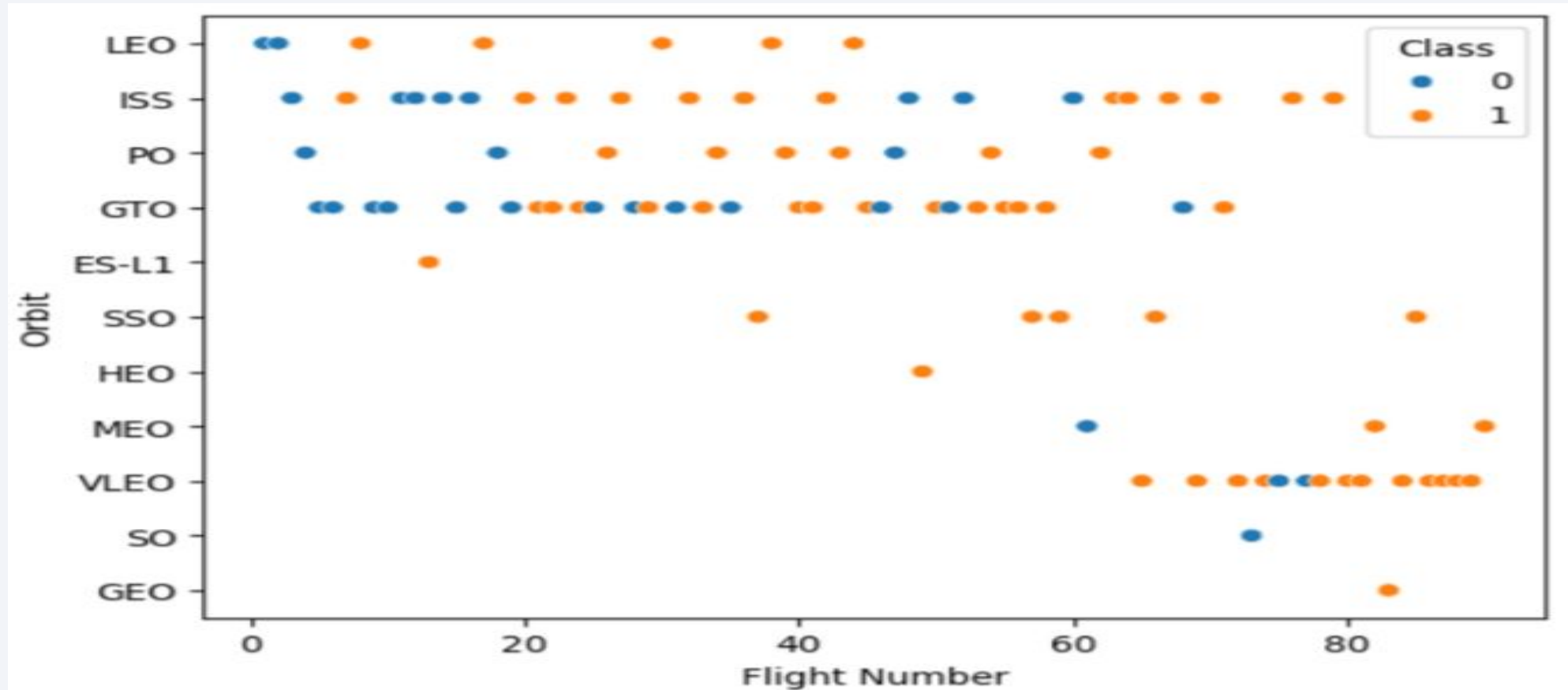- VAFB-SLC has no rockets launched for heavypayload mass (greater than 10000 kg).

# Success Rate vs. Orbit Type

**Explanation**:

- Orbits with 100% success rate:
  - ES-L1
  - GEO
  - HEO
  - SSO
- Orbits with 0% success rate:
  - SO
- Orbits with success rate between 50% and 85%:
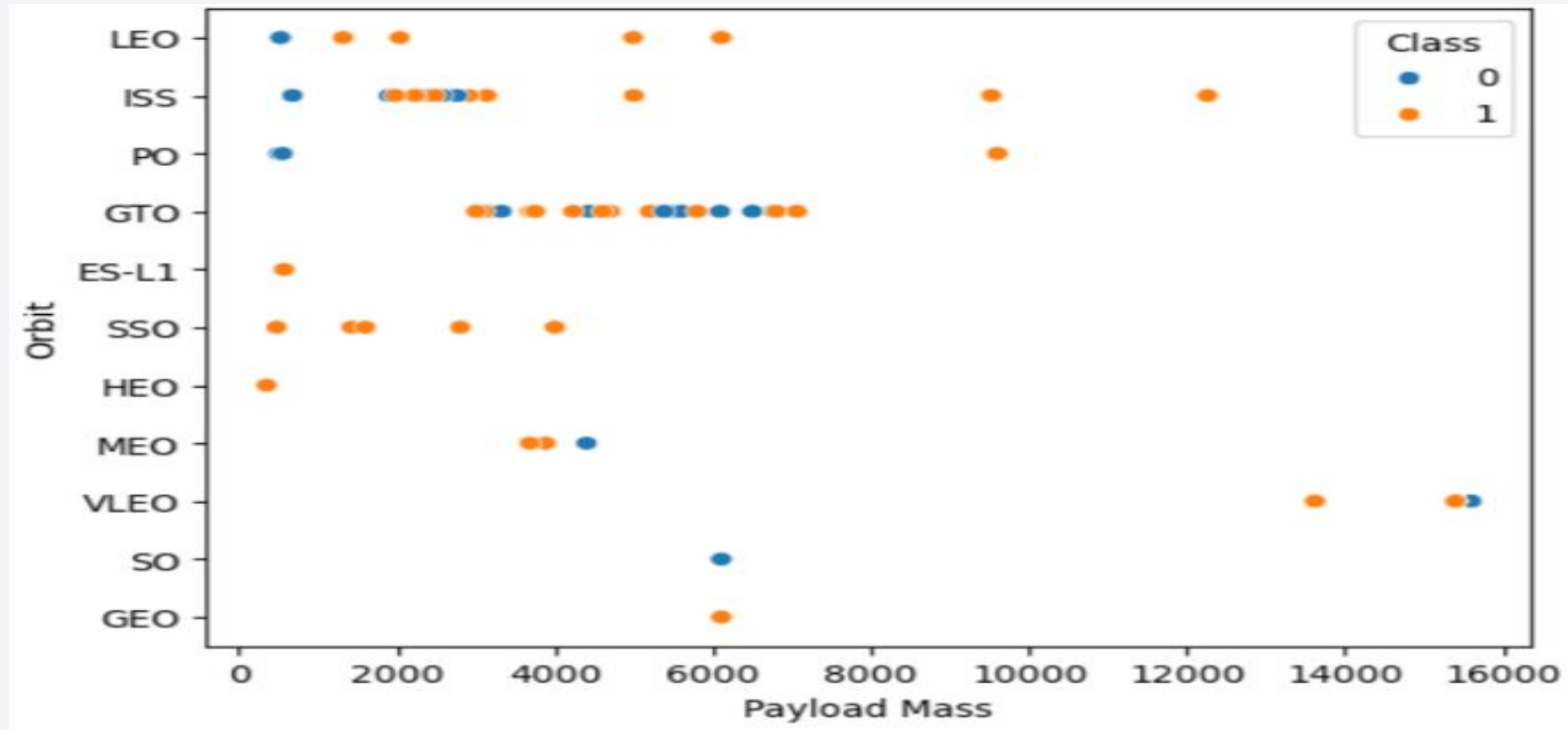  - GTO
  - ISS
  - LEO
  - MEO
  - PO

# Flight Number vs. Orbit Type



- In the LEO orbit, success appears related to the number of flights; on the other hand, there seems to be no relationship between flight number and success in GTO orbit.
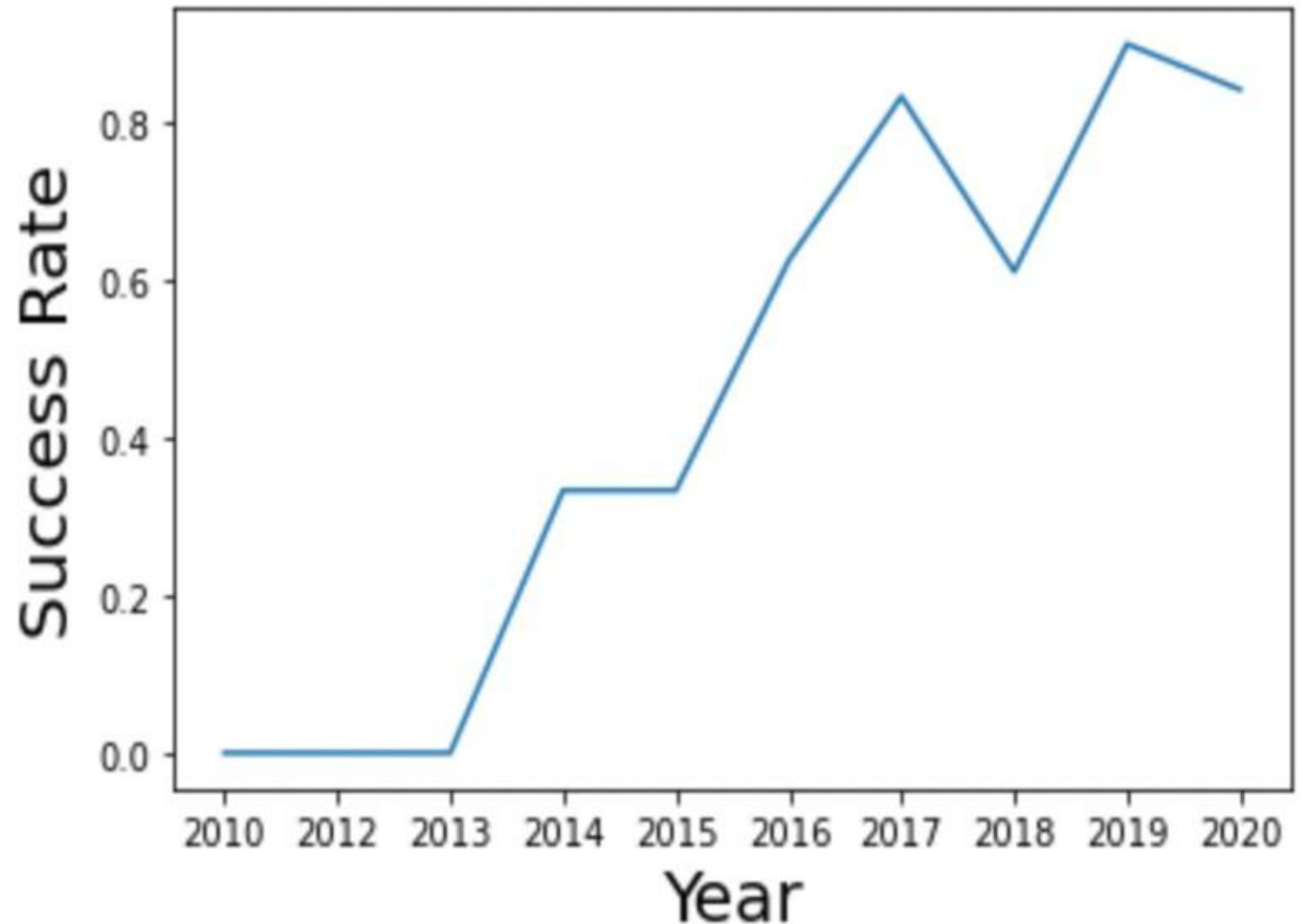
# Payload Mass vs. Orbit Type



- Heavy payloads have a negative influence on GTO orbits and positive on Polar, LEO, and ISS orbits.

# Launch Success Yearly Trend

**Explanation:**

- The success rate since 2013 kept increasing till 2020.

# EDA with SQL

# All Launch Site Names

```
In [4]: %sql select distinct launch_site from SPACEXDATASET;

         * ibm_db_sa://wzf08322:***@0c77d6f2-5da9-48a9-81f8-86b520b87518.bs2io90108kqb1od8lcg.databases.appdomain.cloud:31198/bludb
        Done.

Out[4]:
```

| launch_site |
| --- |
| CCAFS LC-40 |
| CCAFS SLC-40 |
| KSC LC-39A |
| VAFB SLC-4E |

**Explanation**:

- Displaying the names of the unique launch sites in the space mission.

# Launch Site Names Begin with 'CCA'

| 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 | 7:44:00 | F9 v1.0 B0005 | CCAFS LC-40 | Dragon demo flight C2 | 525 | LEO (ISS) | NASA (COTS) | Success | No attempt |
| 2012-10-08 | 0:35:00 | 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 |

**Explanation:**

- Displaying 5 records where launch sites begin with the string 'CCA'.

27

# Total Payload Mass

```
%%sql
select SUM(PAYLOAD_MASS__KG_) as Total_Payload_Mass from Spacextable where Customer like "%CRS%"
```

 * sqlite:///my_data1.db
Done.

**Total_Payload_Mass**

48213

**Explanation:**

- Displaying the total payload mass carried by boosters launched by NASA (CRS).

# Average Payload Mass by F9 v1.1

```sql
%%sql
select Round(AVG(PAYLOAD_MASS__KG_), 3) as Average_Payload_Mass from Spacextable where Booster_Version like "F9 v1.1%"
```

* sqlite:///my_data1.db
Done.

| Average_Payload_Mass |
| --- |
| 2534.667 |

**Explanation:**

- Displaying average payload mass carried by booster version F9 v1.1.

# First Successful Ground Landing Date

```
%%sql
Select Min(Date) as Min_Date, Landing_Outcome from Spacextable where Landing_Outcome like "Success (ground pad)";
```

* sqlite:///my_data1.db
Done.

| Min_Date | Landing_Outcome |
|----------|-----------------|
| 2015-12-22 | Success (ground pad) |

**Explanation**:

- Listing the date when the first successful landing outcome in ground pad was achieved.

# Successful Drone Ship Landing with Payload between 4000 and 6000

| Booster_Version | Payload | PAYLOAD_MASS__KG_ | Landing_Outcome |
|---|---|---|---|
| F9 FT B1022 | JCSAT-14 | 4696 | Success (drone ship) |
| F9 FT B1026 | JCSAT-16 | 4600 | Success (drone ship) |
| F9 FT B1021.2 | SES-10 | 5300 | Success (drone ship) |
| F9 FT B1031.2 | SES-11 / EchoStar 105 | 5200 | Success (drone ship) |

**Explanation:**

- Listing the names of the boosters which have success in drone ship and have payload mass greater than 4000 but less than 6000.

# Total Number of Successful and Failure Mission Outcomes

| Mission_Outcome | Total_Number |
|---|---|
| Failure (in flight) | 1 |
| Success | 98 |
| Success | 1 |
| Success (payload status unclear) | 1 |

**Explanation:**

- Listing the total number of successful and failure mission outcomes.

# Boosters Carried Maximum Payload

| 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 |
| F9 B5 B1051.4 | 15600 |
| F9 B5 B1049.5 | 15600 |
| F9 B5 B1060.2 | 15600 |
| F9 B5 B1058.3 | 15600 |
| F9 B5 B1051.6 | 15600 |
| F9 B5 B1060.3 | 15600 |
| F9 B5 B1049.7 | 15600 |

**Explanation:**

- Listing the names of the booster versions which have carried the maximum payload mass.

# 2015 Launch Records

| Month | Date | Booster_Version | Launch_Site | Landing_Outcome |
|---|---|---|---|---|
| 01 | 2015-01-10 | F9 v1.1 B1012 | CCAFS LC-40 | Failure (drone ship) |
| 04 | 2015-04-14 | F9 v1.1 B1015 | CCAFS LC-40 | Failure (drone ship) |

**Explanation:**

- Listing the failed landing outcomes in drone ship, their booster versions and launch site names for the months in year 2015.

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

| Landing_Outcome | Count_outcome | Date |
|---|---|---|
| No attempt | 10 | 2012-05-22 |
| Success (drone ship) | 5 | 2016-04-08 |
| Failure (drone ship) | 5 | 2015-01-10 |
| Success (ground pad) | 3 | 2015-12-22 |
| Controlled (ocean) | 3 | 2014-04-18 |
| Uncontrolled (ocean) | 2 | 2013-09-29 |
| Failure (parachute) | 2 | 2010-06-04 |
| Precluded (drone ship) | 1 | 2015-06-28 |

**Explanation:**

- 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.
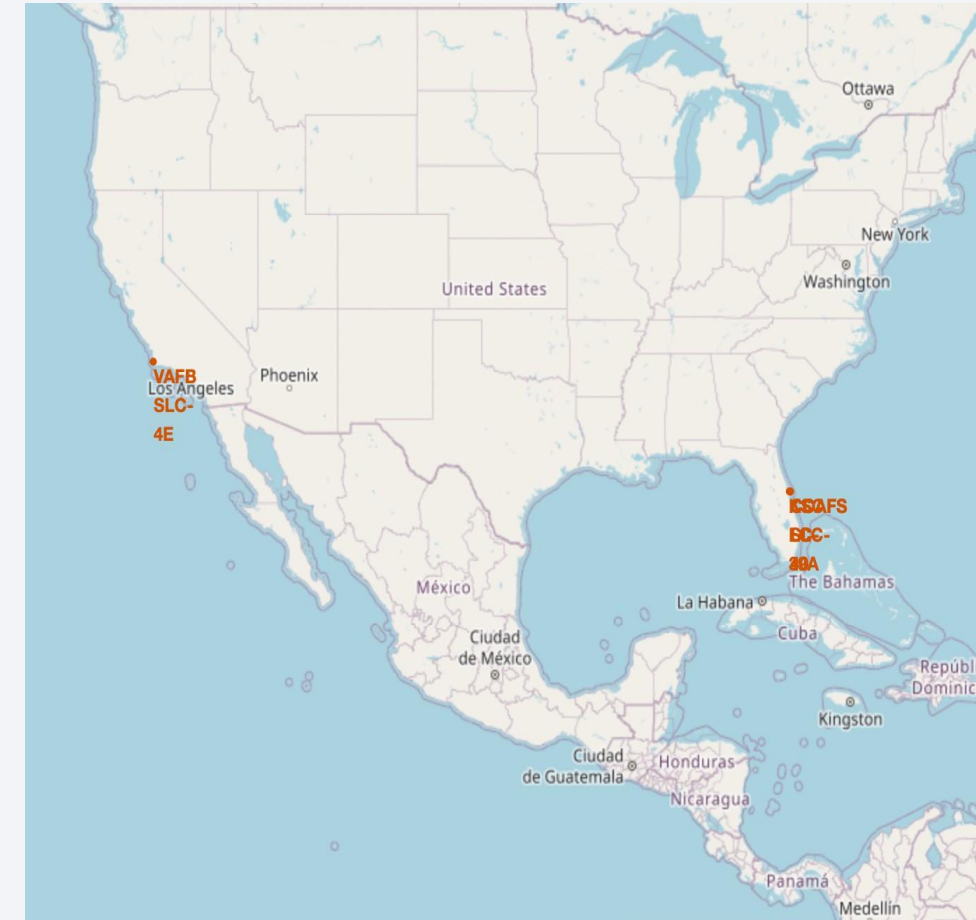
# Launch Sites Proximities Analysis

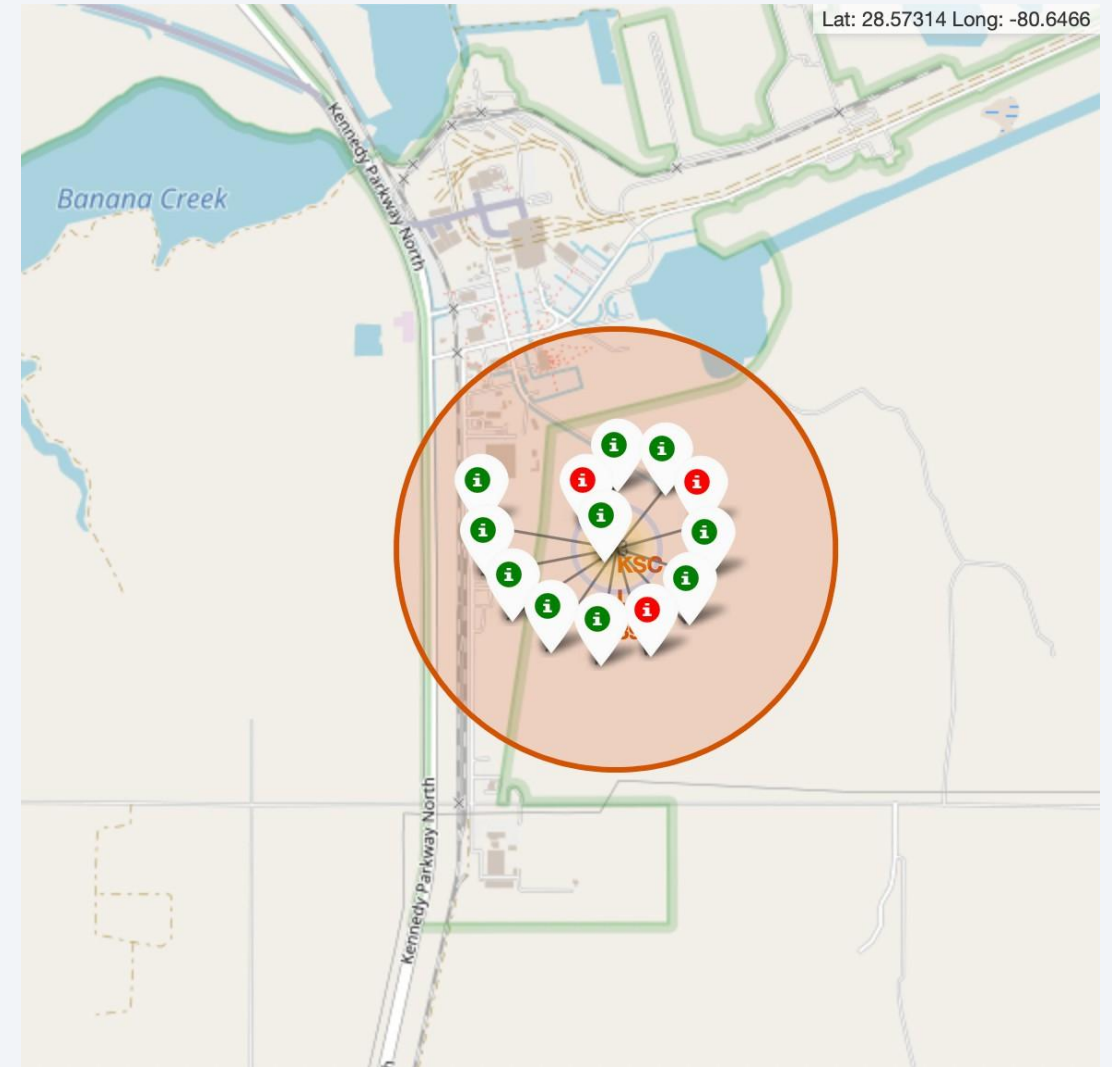# All launch sites' location markers on a Global Map

**Explanation:**

- Most of Launch sites are in proximity to the Equator line. The land is moving faster at the equator than any other place on the surface of the Earth. Anything on the surface of the Earth at the equator is already moving at 1670 km/hour. If a ship is launched from the equator it goes up into space, and it is also moving around the Earth at the same speed it was moving before launching. This is because of inertia. This speed will help the spacecraft keep up a good enough speed to stay in orbit.

- All launch sites are in very close proximity to the coast, while launching rockets towards the ocean it minimises the risk of having any debris dropping or exploding near people.

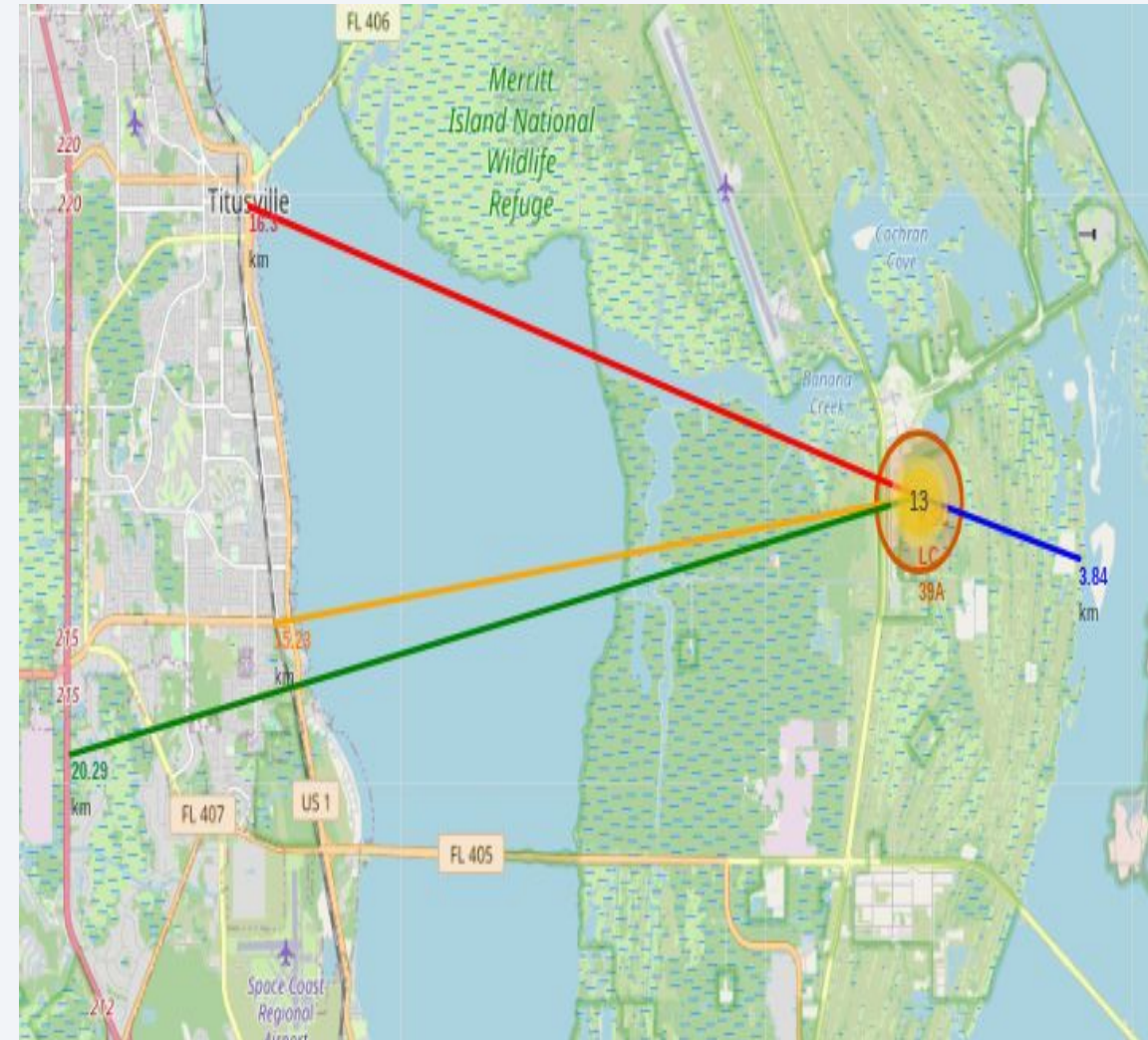# Clustered color labeled launch records on the map

**Explanation:**

- From the colour-labeled markers we should be able to easily identify which launch sites have relatively high success rates.

  - **Green Marker** = Successful Launch
  - **Red Marker** = Failed Launch

- Launch Site KSC LC-39A has a very high Success Rate.

# Distance from the launch site KSC LC-39A to its proximites

**Explanation:**

- From the visual analysis of the launch site KSC LC-39A we can clearly see that it is:
  - relative close to railway (15.23 km)
  - relative close to highway (20.29 km)
  - relative close to coastline (3.84 km)

- Also the launch site KSC LC-39A is relative close to its closest city Titusville (16.32 km).

- Failed rocket with its high speed can cover distances like 15-20 km in few seconds. It could be potentially dangerous to populated areas.

Section 4

# Build a Dashboard
# with Plotly Dash

# Launch success count for all sites

Total Success Launches by Site



Legend:
- KSC LC-39A
- CCAFS SLC-40
- VAFB SLC-4E
- CCAFS LC-40

Pie chart values: 41.2%, 23%, 21.4%, 14.4%

**Explanation:**

- The chart clearly shows that from all the sites, KSC LC-39A has the most successful launches with 41.2%

# Launch site with highest launch success ratio

Total Success Launches for Site KSC LC-39A



**Explanation:**

- KSC LC-39A has the highest launch success rate (76.9%) with 10 successfu and only 3 failed landings.

# Payload Mass vs. Launch Outcome for all sites

**Explanation:**

- The charts show that payloads between 2000 and 5500 kg have the highest success rate.

Section 5

# Predictive Analysis (Classification)

# Classification Accuracy

**Explanation:**

- Based on the scores of the Test Set, we can not confirm which method performs best.
- Same Test Set scores may be due to the small test sample size (18 samples). Therefore, we tested all methods based on the whole Dataset.
- The scores of the whole Dataset confirm that the best model is the Support Vector Machine. This model has not only higher scores, but also the highest accuracy (88%)

Scores and Accuracy of the Test Set

|  | LogReg | SVM | Tree | KNN |
|---|---|---|---|---|
| Jaccard_Score | 0.800000 | 0.800000 | 0.800000 | 0.800000 |
| F1_Score | 0.888889 | 0.888889 | 0.888889 | 0.888889 |
| Accuracy | 0.833333 | 0.833333 | 0.833333 | 0.833333 |

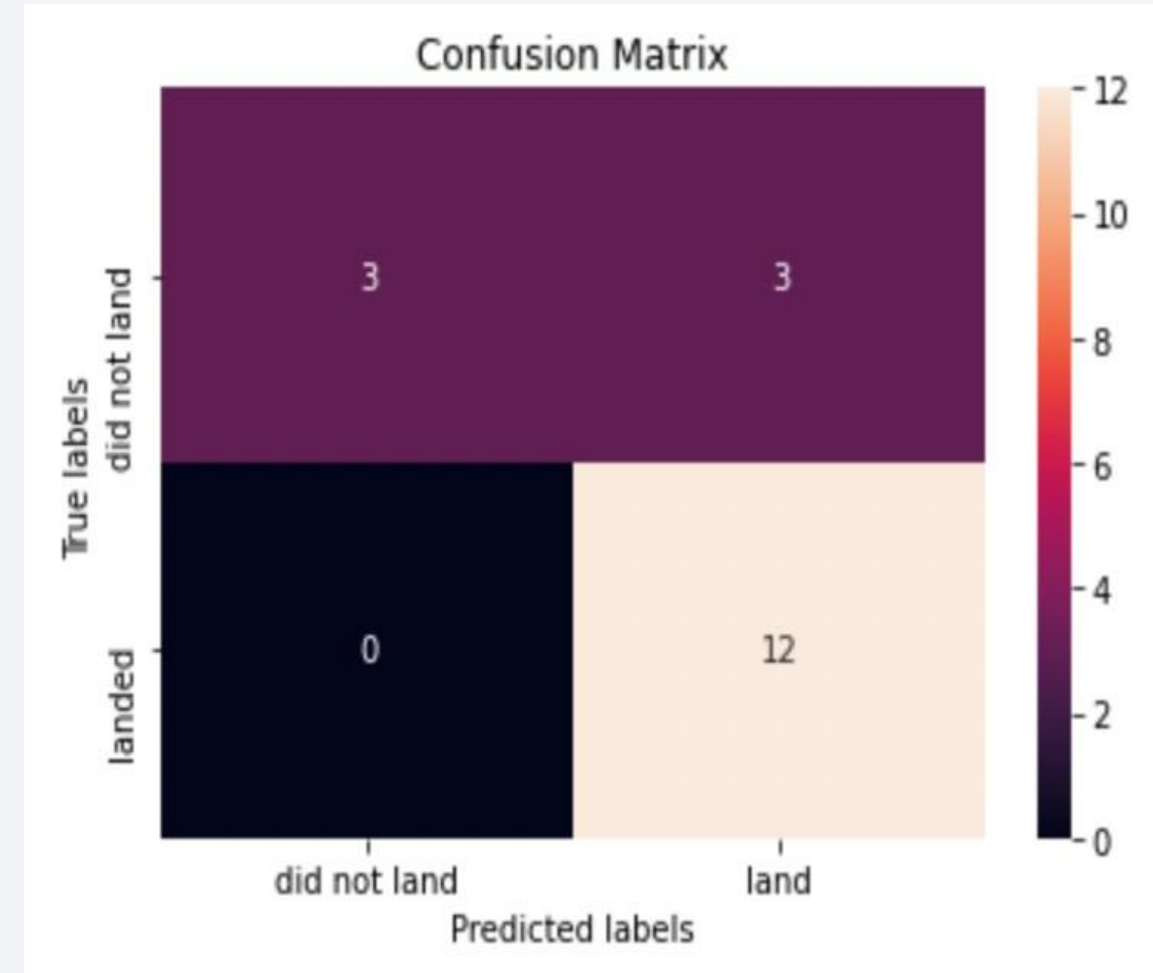Scores and Accuracy of the Entire Data Set

|  | LogReg | SVM | Tree | KNN |
|---|---|---|---|---|
| Jaccard_Score | 0.833333 | 0.845070 | 0.819444 | 0.819444 |
| F1_Score | 0.909091 | 0.916031 | 0.900763 | 0.900763 |
| Accuracy | 0.866667 | 0.877778 | 0.855556 | 0.855556 |

# Confusion Matrix

**Explanation**:

● Examining the confusion matrix, we see that logistic regression can distinguish between the different classes. We see that the major problem is false positives.



Confusion Matrix

|  | Predicted Class | | |
|---|---|---|---|
|  | **Positive** | **Negative** |  |
| **Positive** | True Positive (TP) | False Negative (FN) **Type II Error** | **Sensitivity** $\frac{TP}{(TP+FN)}$ |
| **Negative** | False Positive (FP) **Type I Error** | True Negative (TN) | **Specificity** $\frac{TN}{(TN+FP)}$ |
|  | **Precision** $\frac{TP}{(TP+FP)}$ | **Negative Predictive Value** $\frac{TN}{(TN+FN)}$ | **Accuracy** $\frac{TP+TN}{(TP+TN+FP+FN)}$ |

Actual Class

# Conclusions

- Support Vector Machine (SVM) is the best algorithm for this dataset.

- Launches with a low payload mass show better results than launches with a larger payload mass.

- Most of launch sites are in proximity to the Equator line and all the sites are in very close proximity to the coast.

- The success rate of launches increases over the years.

- KSC LC-39A has the highest success rate of the launches from all the sites.

- Orbits ES-L1, GEO, HEO and SSO have 100% success rate.

# Appendix

**Considerable mentions:**

- Instructors

- Coursera: IBM Professional Data Science Course

- IBM

Thank you!