

IBM Data Science Professional Certificate

Final Capstone Project



ROHAN RAO
31/07/2025

Outline

Executive Summary

Introduction

Methodology

Results

Conclusion

Appendix



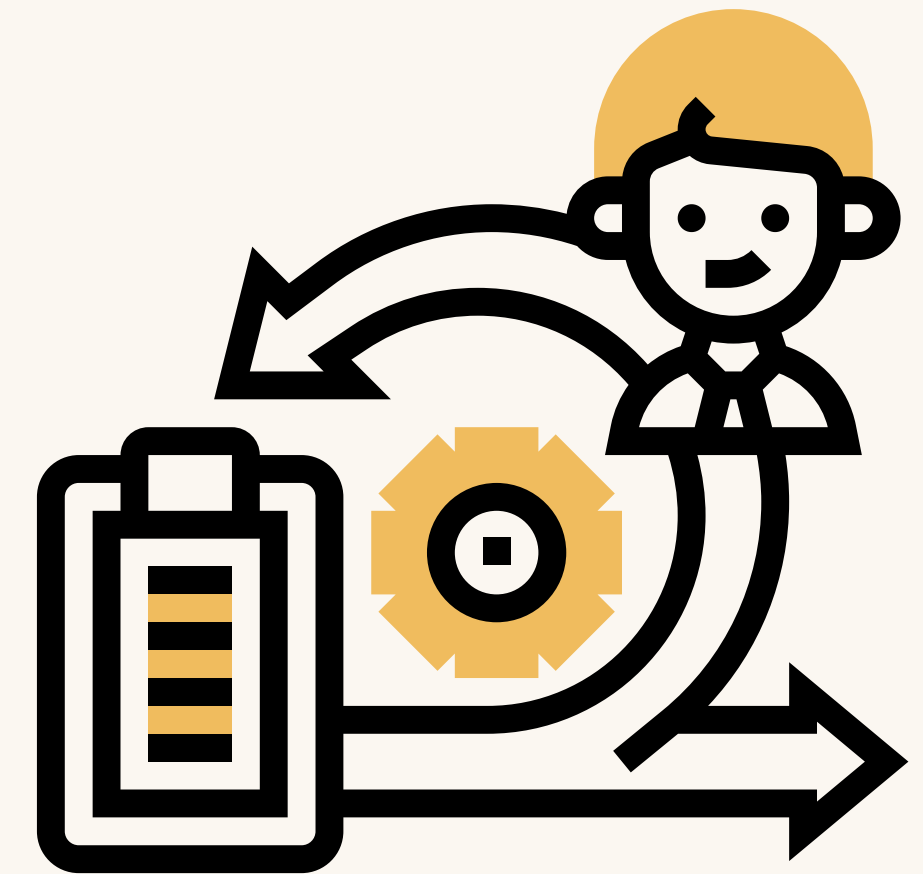
Executive Summary

Summary of methodologies

- Data collection
- Data wrangling
- Exploratory Data Analysis with Data Visualization
- Exploratory Data Analysis with SQL
- Building an interactive map with Folium
- Building a Dashboard with Plotly Dash
- Predictive analysis (Classification)

Summary of all results

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



Introduction

The Falcon 9 rocket, developed by SpaceX, represents a major breakthrough in reusable rocket technology, significantly reducing the cost of space launches by successfully landing its first stage for refurbishment and reuse.

This project aims to analyze historical launch data to predict the success of Falcon 9 first-stage landings using machine learning techniques. By leveraging various launch parameters such as rocket version, payload mass, launch site, and landing type, the objective is to build predictive models that can assess the likelihood of a successful landing. These insights are critical for optimizing launch strategies and improving cost efficiency in the commercial spaceflight industry.

This capstone project demonstrates the practical application of data science and machine learning to solve real-world aerospace challenges.

Data Collection

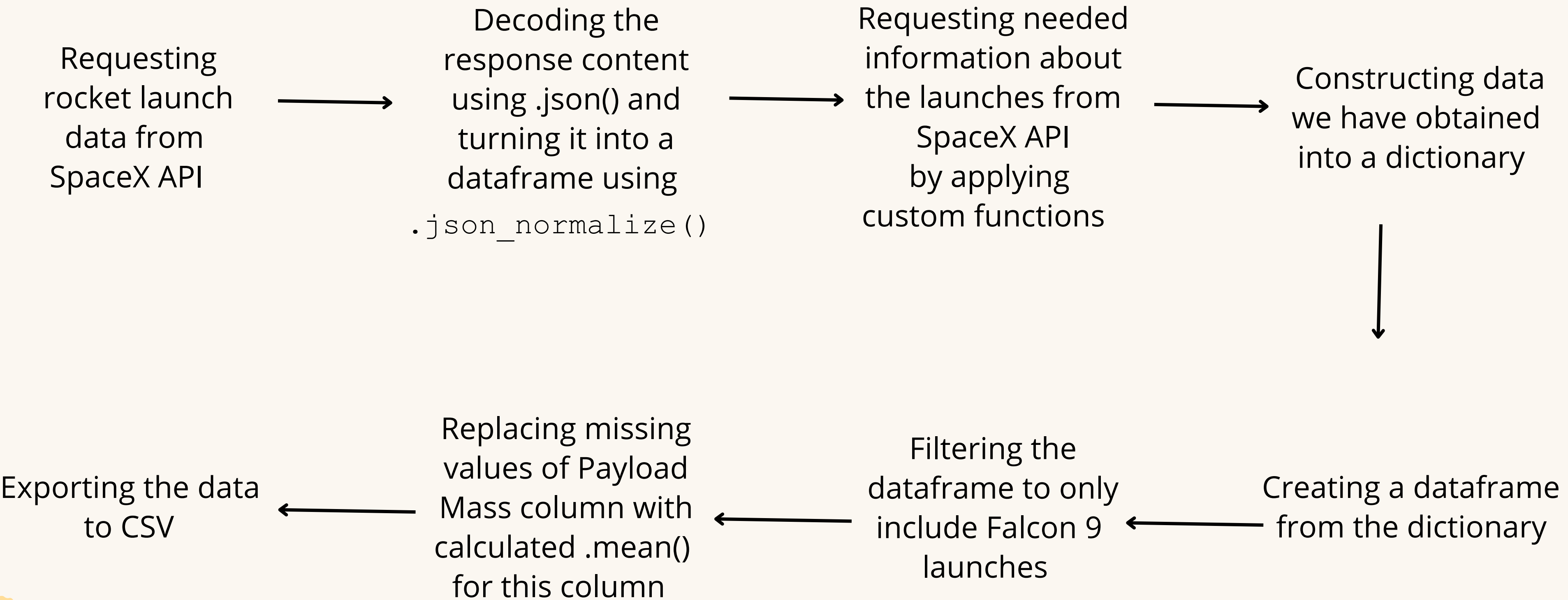
Data collection process involved a combination of API requests from SpaceX REST API and Web Scraping data from a table in SpaceX's Wikipedia entry.

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

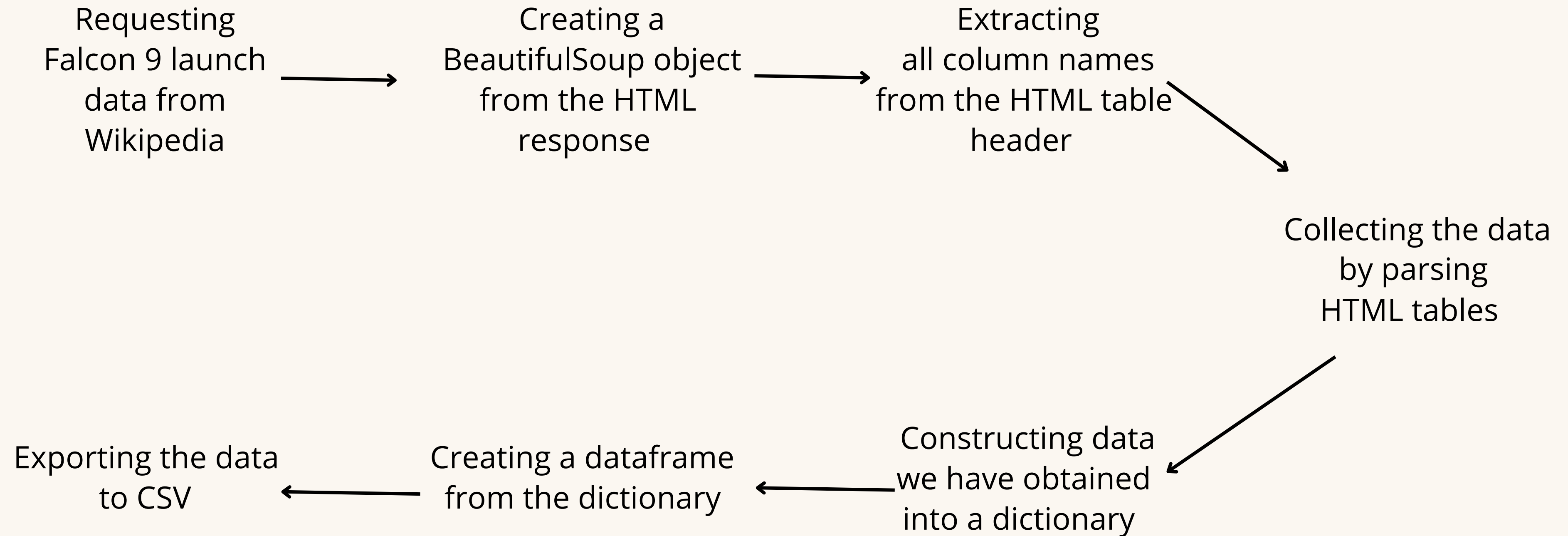
Data Columns are obtained by using SpaceX REST API: FlightNumber, Date, BoosterVersion, PayloadMass, Orbit, LaunchSite, Outcome, Flights, GridFins, Reused, Legs, LandingPad, Block, ReusedCount, Serial, Longitude, Latitude

Data Columns are obtained by using Wikipedia Web Scraping: Flight No., Launch site, Payload, PayloadMass, Orbit, Customer, Launch outcome, Version Booster, Booster landing, Date, Tim

Data Collection - SpaceX Rest API



Data Collection - Web Scrapping



Data Wrangling

Identify which columns are numerical and categorical:

```
In [5]: df.dtypes
```

```
Out[5]: FlightNumber      int64
Date                    object
BoosterVersion          object
PayloadMass             float64
Orbit                   object
LaunchSite              object
Outcome                 object
Flights                 int64
GridFins                bool
Reused                  bool
Legs                    bool
LandingPad              object
Block                   float64
ReusedCount             int64
Serial                  object
Longitude               float64
Latitude               float64
dtype: object
```

Use the method `.value_counts()` to determine the number and occurrence of each orbit in the column `Orbit`

```
: orbit_counts = df['Orbit'].value_counts()
print(orbit_counts)
```

```
Orbit
GTO      27
ISS      21
VLEO     14
PO        9
LEO        7
SSO        5
MEO        3
HEO        1
ES-L1     1
SO         1
GEO        1
Name: count, dtype: int64
```

```
landing_outcomes = df['Outcome'].value_counts()
print(landing_outcomes)
```

```
Outcome
True ASDS      41
None None      19
True RTLS      14
False ASDS       6
True Ocean       5
False Ocean      2
None ASDS        2
False RTLS       1
Name: count, dtype: int64
```


EDA with data visualization

To gain a comprehensive understanding of the Falcon 9 launch dataset, exploratory data analysis (EDA) was performed to identify patterns, trends, and correlations among key variables.

Initial analysis focused on the distribution of landing outcomes across different rocket versions, launch sites, and payload masses. Visualizations played a critical role in this process, bar plots highlighted the frequency of successful landings per launch site, while scatter plots revealed how payload mass impacts landing success. Heatmaps were used to examine correlations between numerical features such as payload mass, orbit type, and success probability. Additionally, pie charts and histograms provided insights into categorical distributions such as launch outcomes and booster versions.

These visualizations helped uncover that certain launch sites and payload ranges had higher landing success rates, laying the groundwork for feature selection in the machine learning phase.

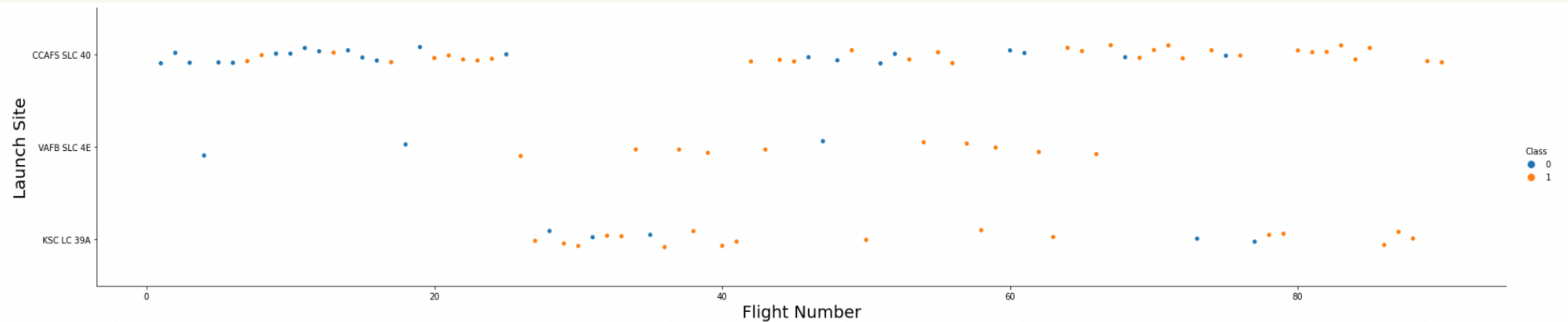
EDA with SQL

To gather preliminary insights from the SpaceX Falcon 9 launch dataset, structured exploratory data analysis was performed using SQL. The following key queries were executed:

1. Display the names of the unique launch sites in the space mission dataset to understand the geographical distribution of Falcon 9 launches.
2. Display 5 records where launch sites begin with the string 'CCA', providing an initial look into launch activity at Cape Canaveral sites.
3. Display the total payload mass carried by boosters launched by NASA (CRS) to evaluate the scale of NASA's missions.
4. Display the average payload mass carried by booster version F9 v1.1, offering insight into the typical capacity handled by earlier booster variants.
5. List the date when the first successful landing outcome on a ground pad was achieved, highlighting SpaceX's milestone in reusability.
6. List the names of the boosters which had success on a drone ship and a payload mass between 4000 and 6000 kg, identifying high-performing missions in a specific payload range.
7. List the total number of successful and failed mission outcomes, to assess overall mission reliability.
8. List the names of the booster versions which have carried the maximum payload mass, revealing the highest-capacity hardware.
9. List the failed landing outcomes on drone ships, their booster versions, and launch site names for 2015, focusing on failures during a crucial developmental phase.
10. Rank the count of landing outcomes (e.g., "Failure (drone ship)" or "Success (ground pad)") between the dates 2010-06-04 and 2017-03-20, in descending order, to understand which outcomes were most frequent during that time frame.

EDA WITH VISUALIZATION

Flight Number vs. Launch Site



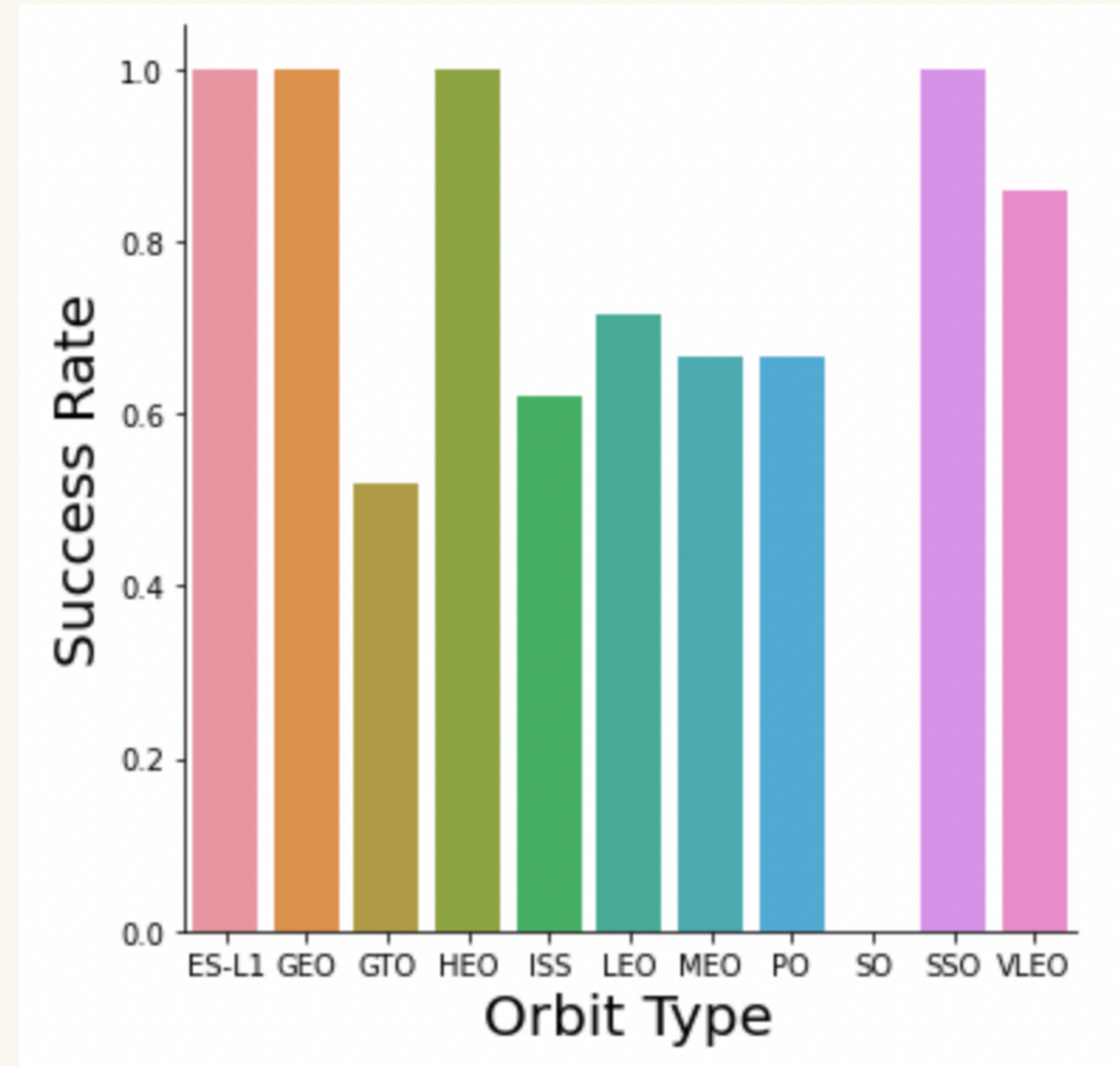
Explanation:

- The earliest flights all failed while the latest flights all succeeded.
- The CCAFS SLC 40 launch site has about a half of all launches.
- 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.

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



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.bs2io90108kqblod8lcg.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`

`In [5]: %sql select * from SPACEXDATASET where launch_site like 'CCA%' limit 5;`

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

`Out[5]:`

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

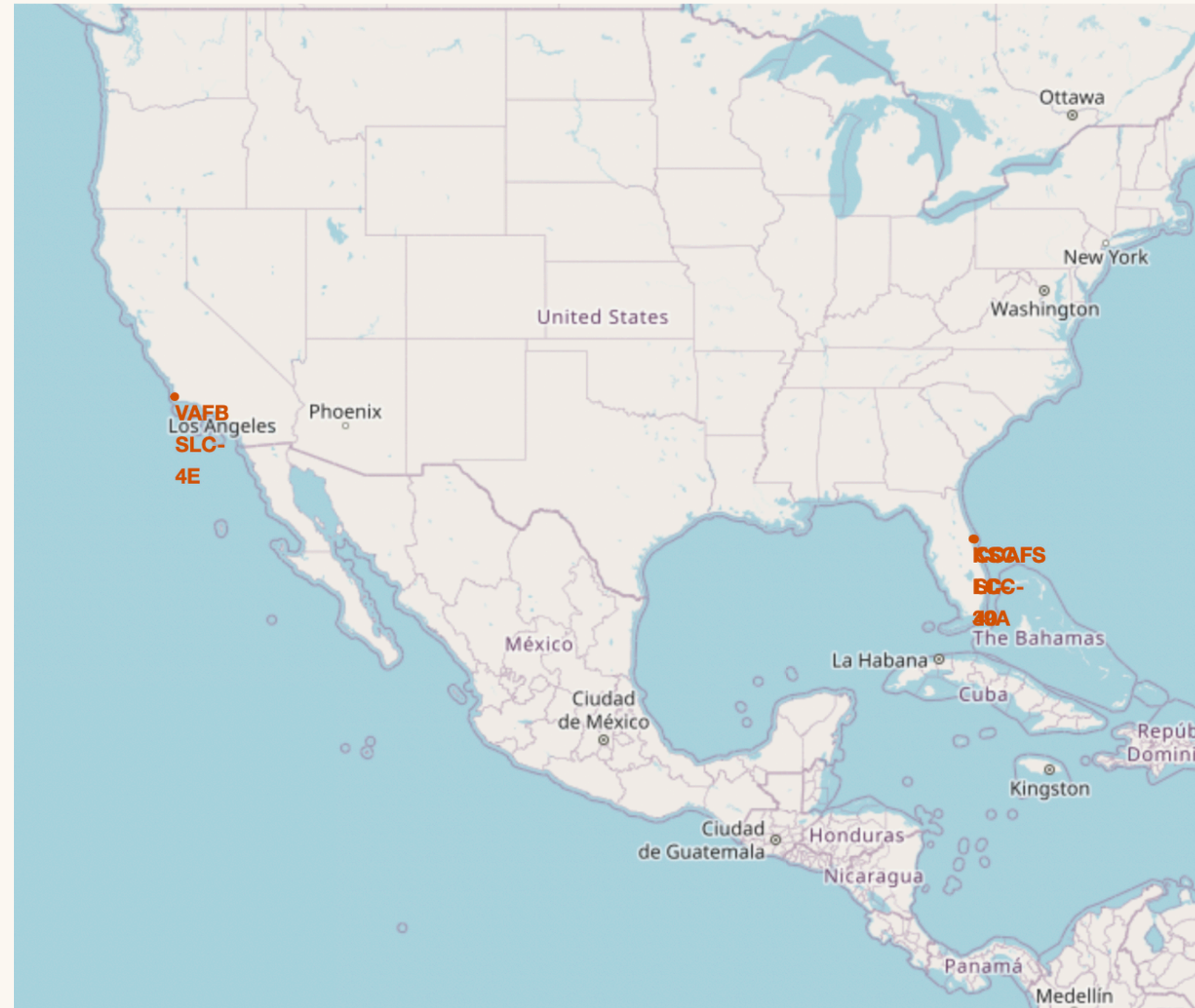
Explanation:

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

INTERACTIVE MAP USING FOLIUM

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.



DASHBOARD USING PLOTLY DASH

Total Success Launches by Site



Explanation:

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

Predictive analysis (classification)

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 Decision Tree Model. This model has not only higher scores, but also the highest accuracy.

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

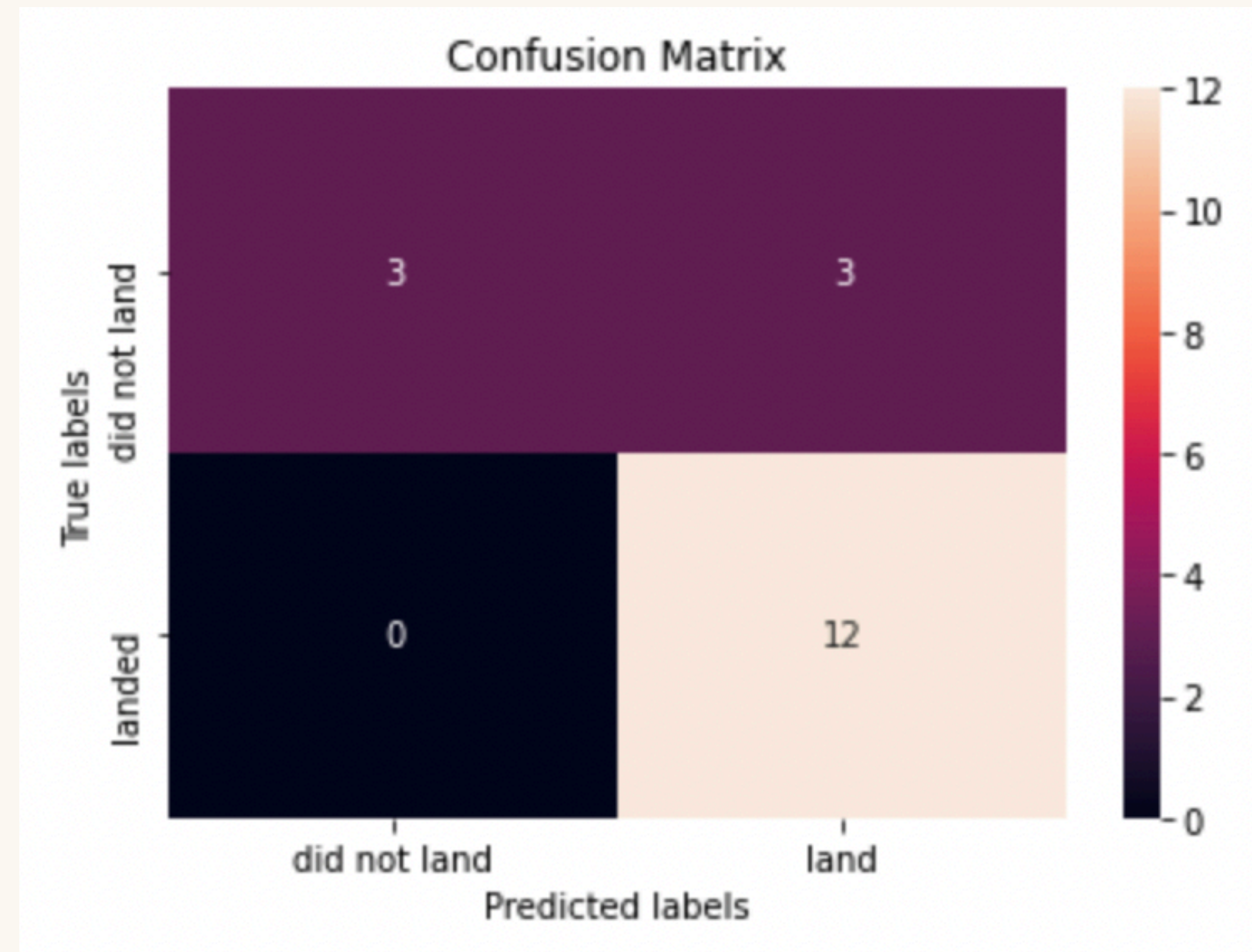
	LogReg	SVM	Tree	KNN
Jaccard_Score	0.833333	0.845070	0.882353	0.819444
F1_Score	0.909091	0.916031	0.937500	0.900763
Accuracy	0.866667	0.877778	0.911111	0.855556

Confusion Matrix

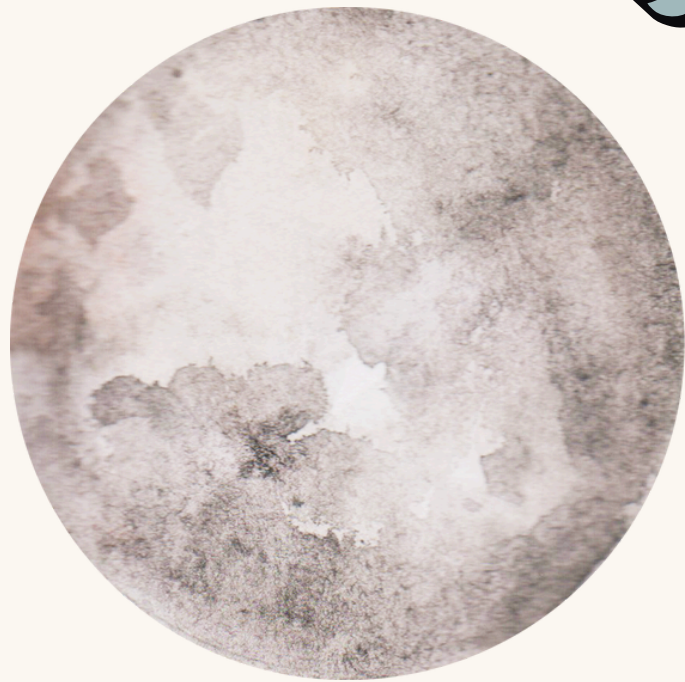
Explanation:

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

		Predicted Values	
		Negative	Positive
Actual Values	Negative	TN	FP
	Positive	FN	TP



Conclusion



- Decision Tree Model 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.