



IBM Developer
SKILLS NETWORK

Winning Space Race with Data Science

<Badis Nahdi>

<2023-02-04>



Outline

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

Executive Summary

- Summary of methodologies
 - Data Collection with API • Data Collection with Web Scraping
 - Data Wrangling
 - Exploratory Data Analysis with SQL and Data Visualization
 - Interactive Visual Analysis with Folium and Dash
 - Machine Learning Prediction
- Summary of all results
 - Exploratory Data Analysis result • Interactive analytics in screenshots
 - Predictive Analytics result (Machine Learning Lab)

Introduction

- Project background and context
 - SpaceX 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. This information can be used if an alternate
 - company wants to bid against SpaceX for a rocket launch.
- Problems you want to find answers
 - What features determine the landing outcomes?
 - What are the relationships between those features?
 - What state of these features determine the best outcome?

Section 1

Methodology

Methodology

Executive Summary

- Data collection methodology:
 - Describe how data was collected
- Perform data wrangling
 - Describe how data was processed
- 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

Data collection is the process of gathering data for use in business decision-making, strategic planning, research and other purposes. It's a crucial part of data analytics applications and research projects: Effective data collection provides the information that's needed to answer questions, analyse business performance or other outcomes, and predict future trends.

- We used the SpaceX REST API to collect the Data.
- We also used Web Scraping to collect Data from Wikipedia.
- Next, we present our data collection process.

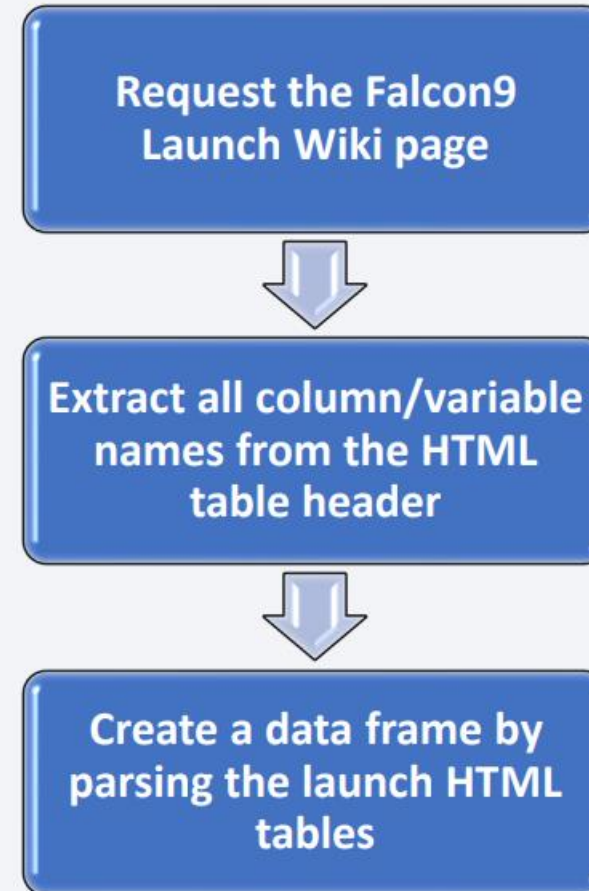
Data Collection – SpaceX API

- Get request for rocket launch data using API
- Use `json_normalize()` method to convert JSON result to Dataframe
- For more details:



Data Collection - Scraping

- Present your web scraping process using key phrases and flowcharts
- Add the GitHub URL of the completed web scraping notebook, as an external reference and peer-review purpose



Data Wrangling

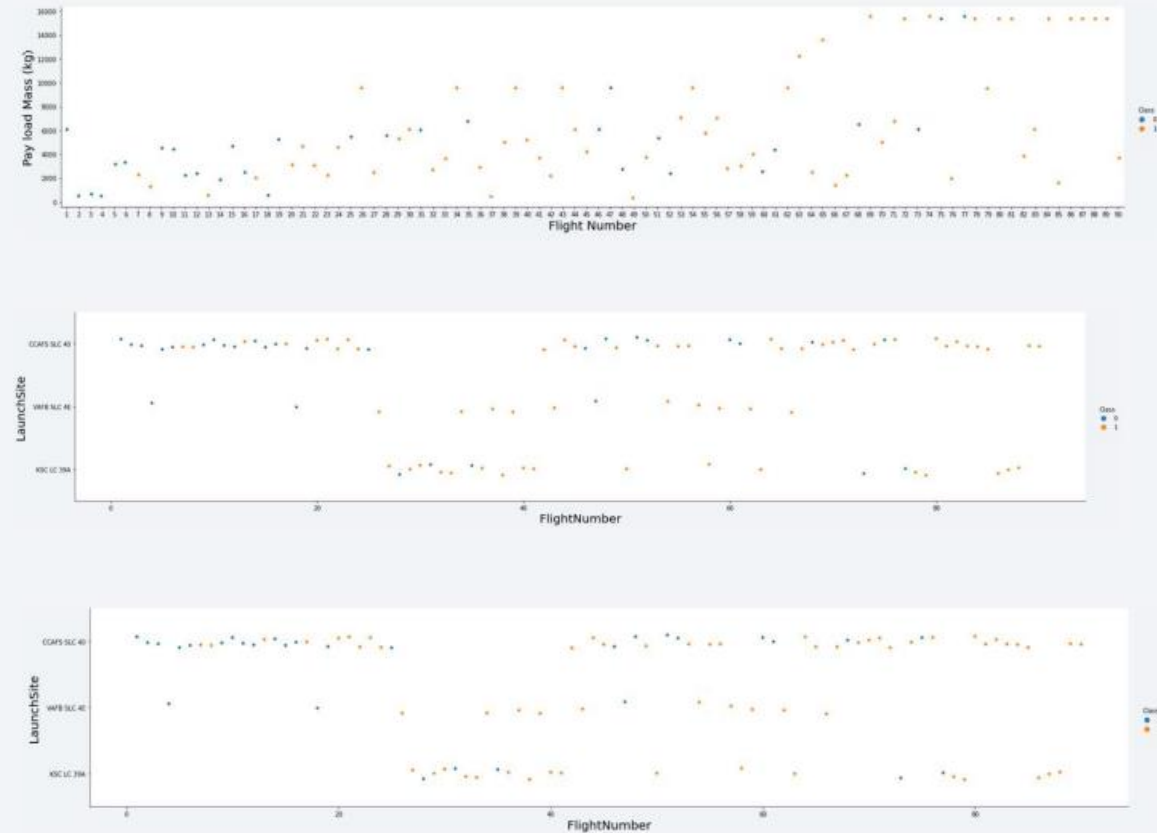
- Data wrangling is the process of cleaning and unifying messy and complex data sets for easy access and analysis. With the amount of data and data sources rapidly growing and expanding, it is getting increasingly essential for large amounts of available data to be organized for analysis.
- We calculate the number of launches at each site, and the number and occurrence of each orbits
- We create landing outcome label from outcome column and export the results to csv.

EDA with Data Visualization

- We used scatter graph to find the relationship between the attributes such as between:

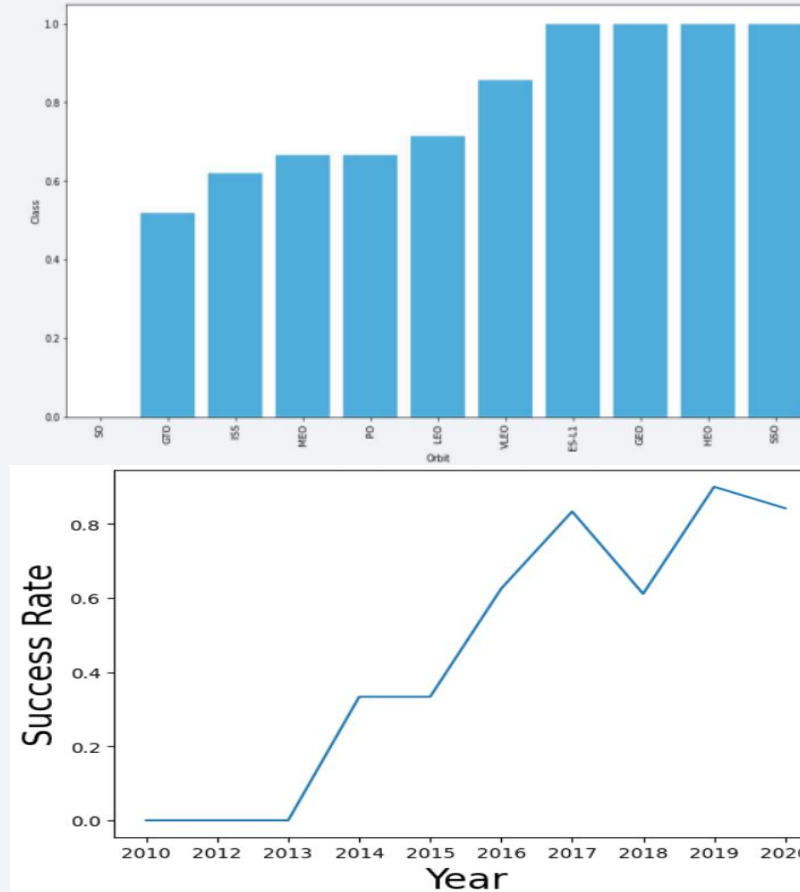
- Payload and Flight Number.
- Flight Number and Launch Site.
- Payload and Launch Site.
- Flight Number and Orbit Type.
- Payload and Orbit Type.

Scatter plots show dependency of attributes on each other. Once a pattern is determined from the graphs. It's very easy to see which factors affecting the most to the success of the landing outcomes.



EDA with Data Visualization

- Bar graphs is one of the easiest way to interpret the relationship between the attributes. In this case, we used a bar chart to determine which orbits have the highest probability of success.
- We used a line graph to show a trends or pattern of the attribute over time which in this case, is used for see the launch success yearly trend.



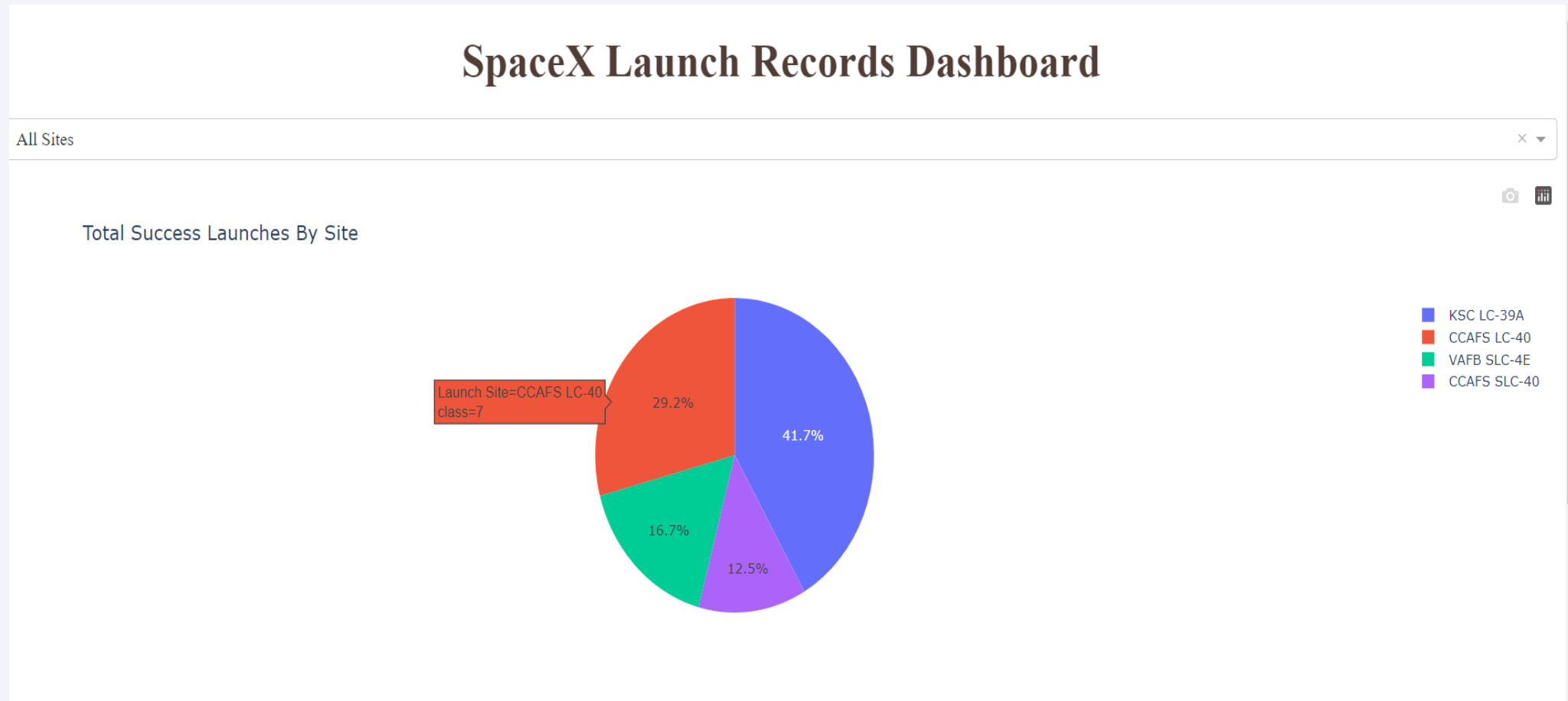
EDA with SQL

- We loaded the SpaceX dataset into an IBM db2 database. We applied EDA with SQL to get insight from the data. We wrote queries to find out for instance:
 - the names of the launch sites.
 - 5 records where launch sites begin with the string 'CCA'.
 - the total payload mass carried by booster launched by NASA (CRS).
 - the average payload mass carried by booster version F9 v1.1.
 - the date when the first successful landing outcome in ground pad was achieved.
 - the names of the boosters which have success in drone ship and have payload mass greater than 4000 but less than 6000.
 - the total number of successful and failure mission outcomes.
 - the names of the booster_versions which have carried the maximum payload mass.
 - the failed landing_outcomes in drone ship, their booster versions, and launch sites names for in year 2015.
 - the count of landing outcomes or success between the date 2010-06-04 and 2017-03-20, in descending order.

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



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

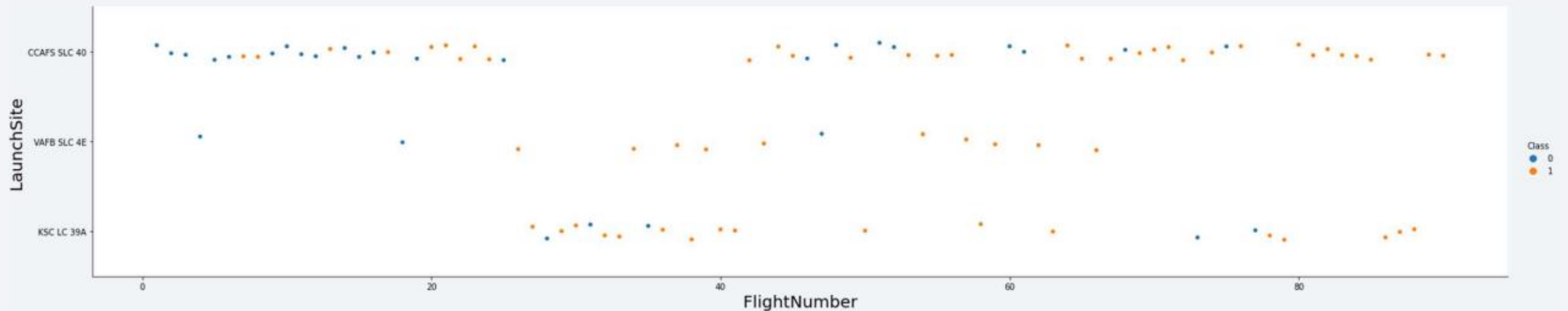
- Exploratory data analysis results
- Interactive analytics demo in screenshots
- Predictive analysis results

The background of the slide is an abstract composition. It features a dark blue base color. Overlaid on this are numerous diagonal streaks in shades of red and cyan. A faint, light blue grid pattern is also visible, particularly in the lower half of the image. The overall effect is dynamic and technological.

Section 2

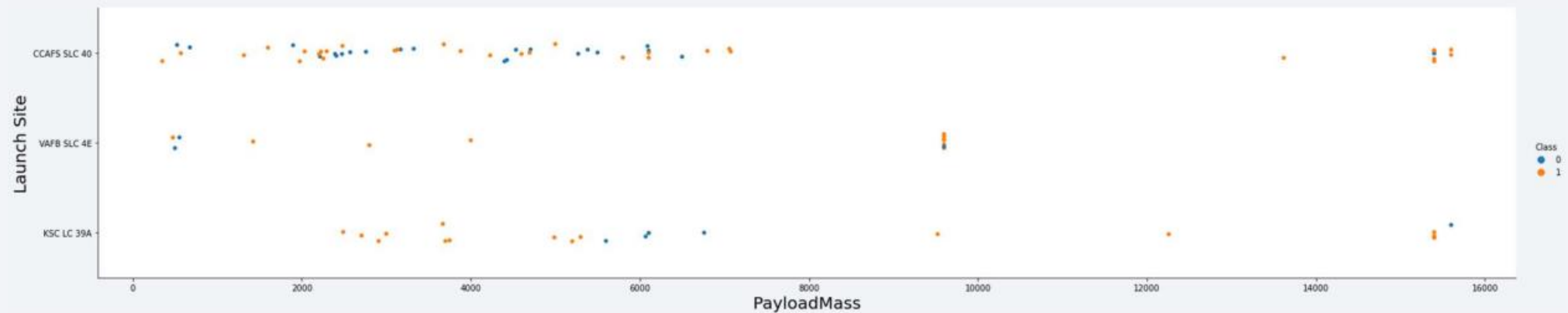
Insights drawn from EDA

Flight Number vs. Launch Site



This plot shows that the larger the flights number of each launch site, the greater the success rate will be.

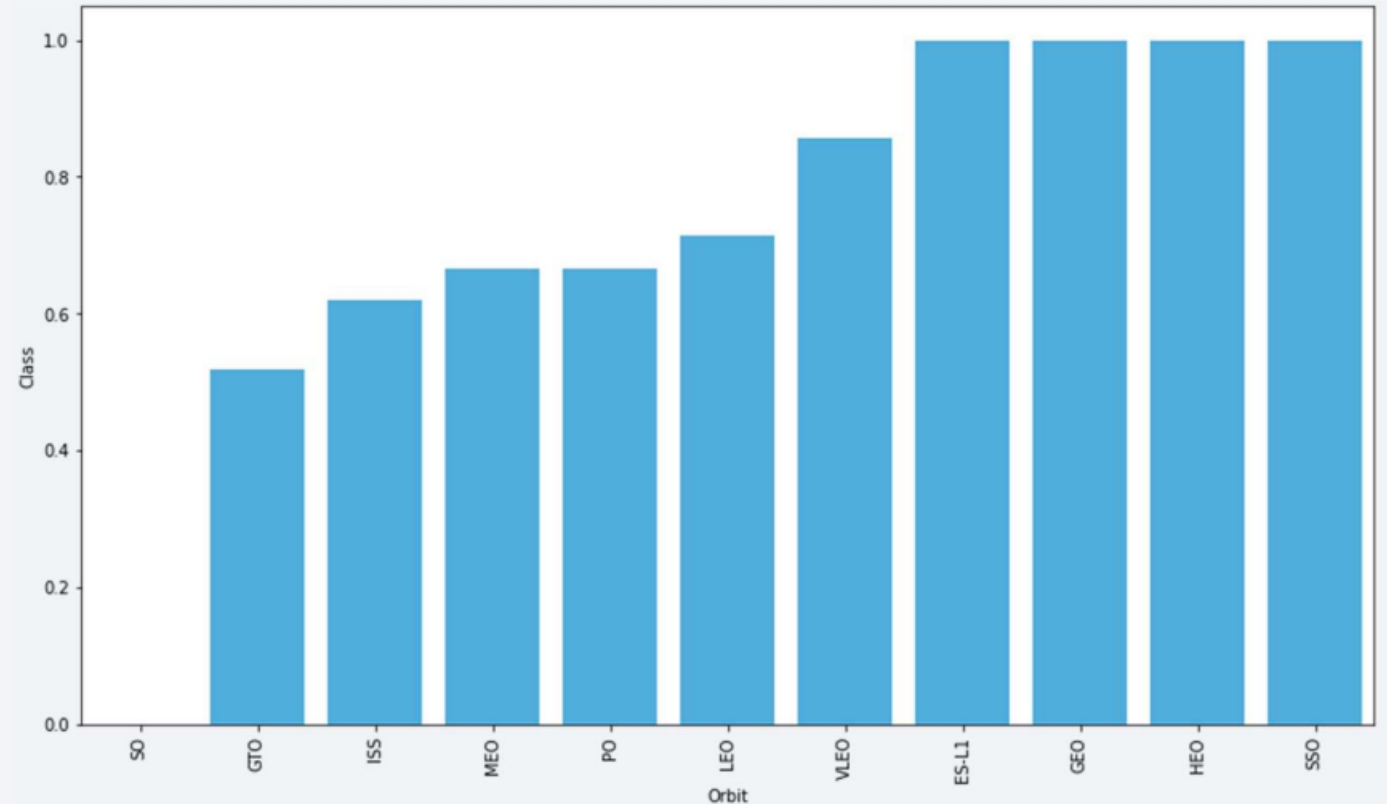
Payload vs. Launch Site



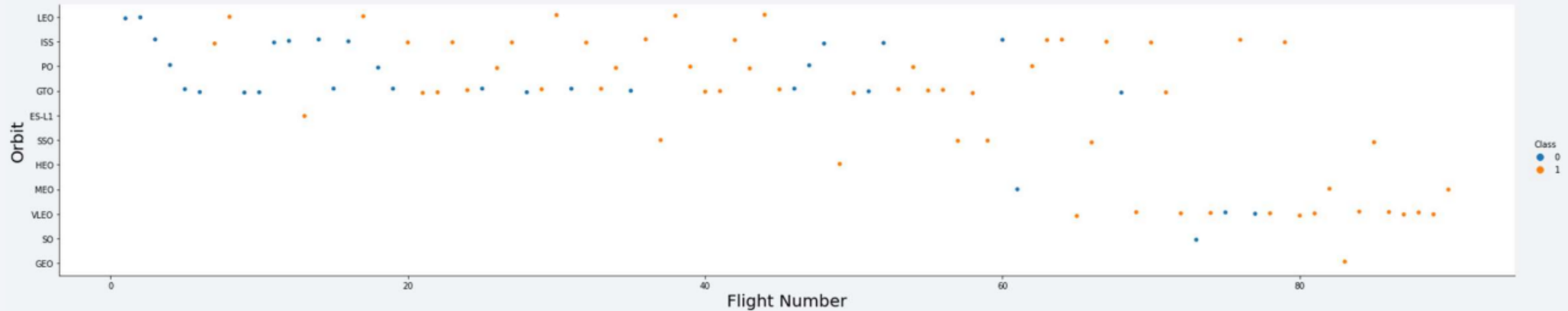
This scatter plot shows once the payload mass is greater than 7000kg, the probability of the success rate will be highly increased. However, there is no clear pattern to say the launch site is dependent to the payload mass for the success rate.

Success Rate vs. Orbit Type

This chart depicted the possibility of the orbits to influences the landing outcomes as some orbits has 100% success rate such as SSO, HEO, GEO AND ES-L1 while SO orbit produced 0% rate of success.



Flight Number vs. Orbit Type



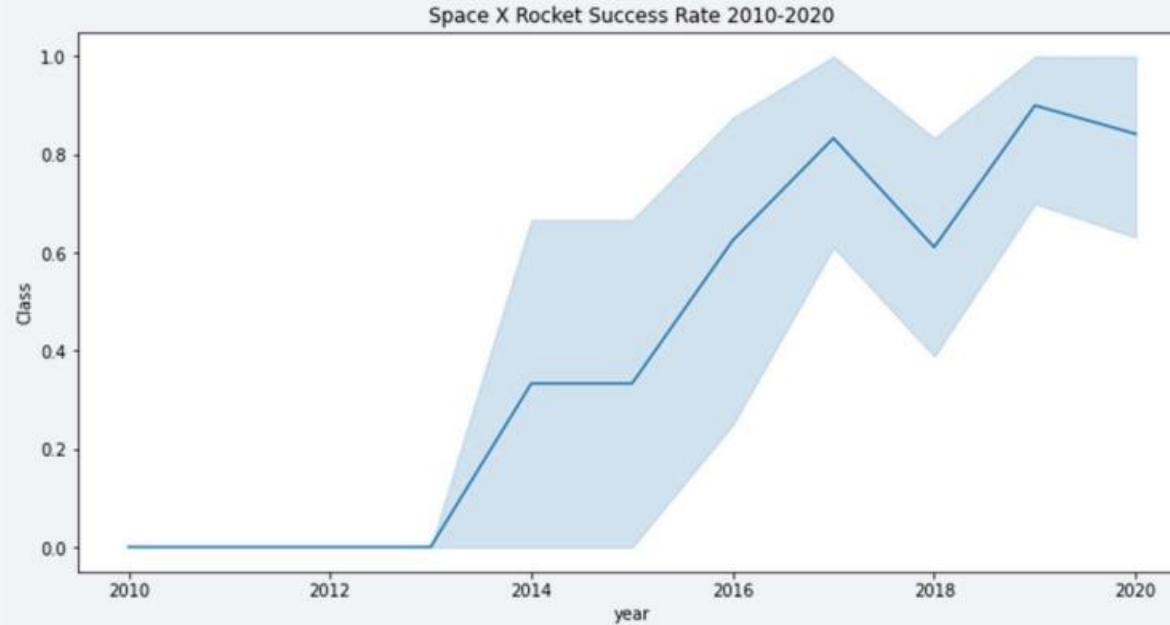
We see that in the LEO orbit the Success appears related to the number of flights; on the other hand, there seems to be no relationship between flight number when in GTO orbit.

Payload vs. Orbit Type



With heavy payloads the successful landing or positive landing rate are more for Polar, LEO and ISS. However, for GTO we cannot distinguish this well as both positive landing rate and negative landing(unsuccesful mission) are both there here.

Launch Success Yearly Trend



you can observe that the success rate since 2013
kept increasing till 2020

All Launch Site Names

```
%sql SELECT DISTINCT "LAUNCH_SITE" FROM SPACEXTBL
```

```
* 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 SPACEXTBL 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
04-06-2010	18:45:00	F9 v1.0 B0003	CCAFS LC-40	Dragon Spacecraft Qualification Unit	0	LEO	SpaceX	Success	Failure (parachute)
08-12-2010	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)
22-05-2012	07:44:00	F9 v1.0 B0005	CCAFS LC-40	Dragon demo flight C2	525	LEO (ISS)	NASA (COTS)	Success	No attempt
08-10-2012	00:35:00	F9 v1.0 B0006	CCAFS LC-40	SpaceX CRS-1	500	LEO (ISS)	NASA (CRS)	Success	No attempt
01-03-2013	15:10:00	F9 v1.0 B0007	CCAFS LC-40	SpaceX CRS-2	677	LEO (ISS)	NASA (CRS)	Success	No attempt

Total Payload Mass

```
%sql SELECT SUM("PAYLOAD_MASS_KG_") FROM SPACEXTBL WHERE "CUSTOMER" = 'NASA (CRS)'  
[10]  
... * sqlite:///my\_data1.db  
Done.  
SUM("PAYLOAD_MASS_KG_")  
45596
```

Average Payload Mass by F9 v1.1

```
[11] %sql SELECT AVG("PAYLOAD_MASS_KG_") FROM SPACEXTBL WHERE "BOOSTER_VERSION" LIKE '%F9 v1.1%'
...  * sqlite:///my\_data1.db
      Done.
</>  AVG("PAYLOAD_MASS_KG_")
      2534.6666666666665
```


First Successful Ground Landing Date

```
[12] %sql SELECT MIN("DATE") FROM SPACEXTBL WHERE "Landing _Outcome" LIKE '%Success%'
...  * sqlite:///my\_data1.db
Done.
</> MIN("DATE")
01-05-2017
```

Successful Drone Ship Landing with Payload between 4000 and 6000

```
%sql SELECT "BOOSTER_VERSION" FROM SPACEXTBL WHERE "LANDING _OUTCOME" = 'Success (drone ship)'
AND "PAYLOAD_MASS_KG_" > 4000 AND "PAYLOAD_MASS_KG_" < 6000;

[13]
... * sqlite:///my_data1.db
Done.

</>
Booster_Version
F9 FT B1022
F9 FT B1026
F9 FT B1021.2
F9 FT B1031.2
```

Total Number of Successful and Failure Mission Outcomes

```
> %sql SELECT (SELECT COUNT("MISSION_OUTCOME") FROM SPACEXTBL WHERE "MISSION_OUTCOME" LIKE '%Success%') AS SUCCESS,  
  (SELECT COUNT("MISSION_OUTCOME") FROM SPACEXTBL WHERE "MISSION_OUTCOME" LIKE '%Failure%') AS FAILURE  
[14]  
.. * sqlite:///my\_data1.db  
Done.  
/>  
SUCCESS FAILURE  
100 1
```

+ Code + Markdown

Boosters Carried Maximum Payload

```
%sql SELECT DISTINCT "BOOSTER_VERSION" FROM SPACEXTBL \
WHERE "PAYLOAD_MASS_KG_" = (SELECT max("PAYLOAD_MASS_KG_") FROM SPACEXTBL)

* 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 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
```

2015 Launch Records

```
%sql SELECT substr("DATE", 4, 2) AS MONTH, "BOOSTER_VERSION", "LAUNCH_SITE" FROM SPACEXTBL\
WHERE "LANDING_OUTCOME" = 'Failure (drone ship)' and substr("DATE",7,4) = '2015'
```

```
* sqlite:///my\_data1.db
```

```
Done.
```

MONTH	Booster_Version	Launch_Site
01	F9 v1.1 B1012	CCAFS LC-40
04	F9 v1.1 B1015	CCAFS LC-40

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

```
%sql SELECT "LANDING _OUTCOME", COUNT("LANDING _OUTCOME") FROM SPACEXTBL\
WHERE "DATE" >= '04-06-2010' and "DATE" <= '20-03-2017' and "LANDING _OUTCOME" LIKE '%Success%'\
GROUP BY "LANDING _OUTCOME" \
ORDER BY COUNT("LANDING _OUTCOME") DESC ;
```

17]

* [sqlite:///my_data1.db](#)

Done.

>

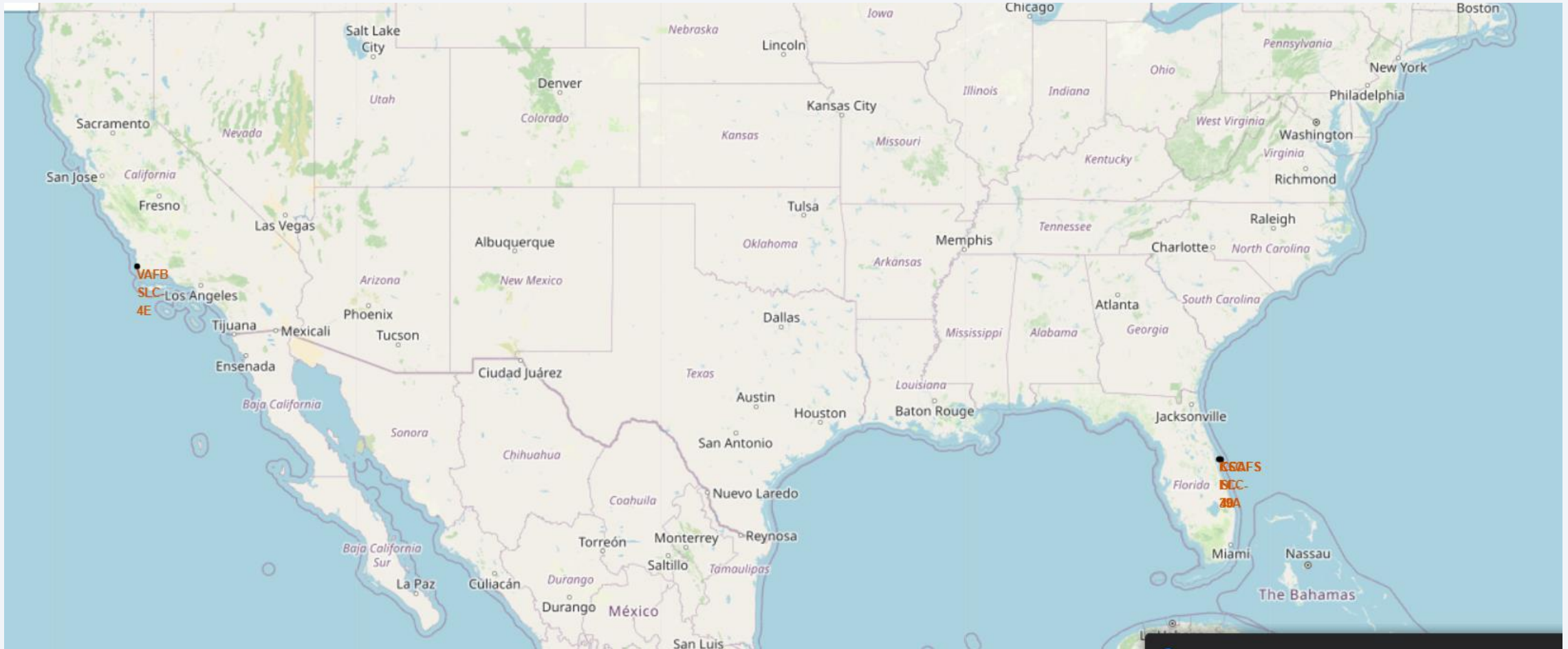
Landing _Outcome	COUNT("LANDING _OUTCOME")
Success	20
Success (drone ship)	8
Success (ground pad)	6

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.

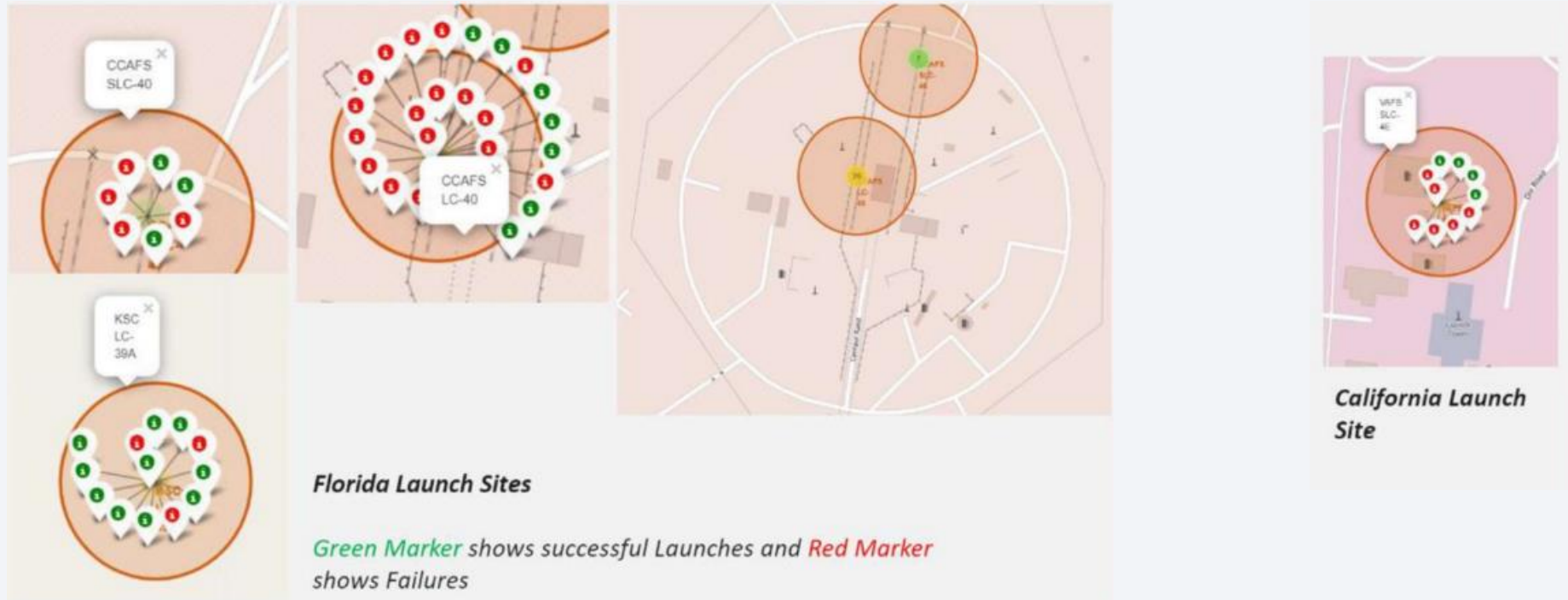
Section 3

Launch Sites Proximities Analysis

Location of all the Launch Sites



Markers showing launch sites with color labels





Section 4

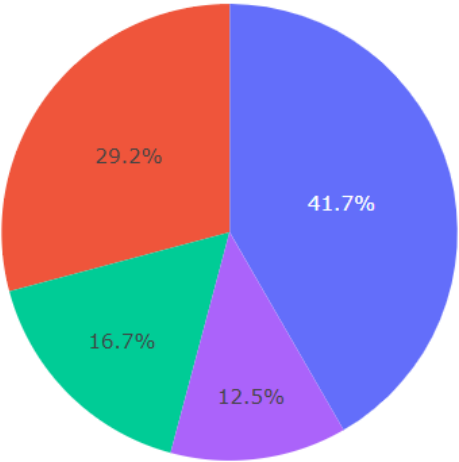
Build a Dashboard with Plotly Dash

SpaceX Launch Records Dashboard

All Sites



Total Success Launches By Site

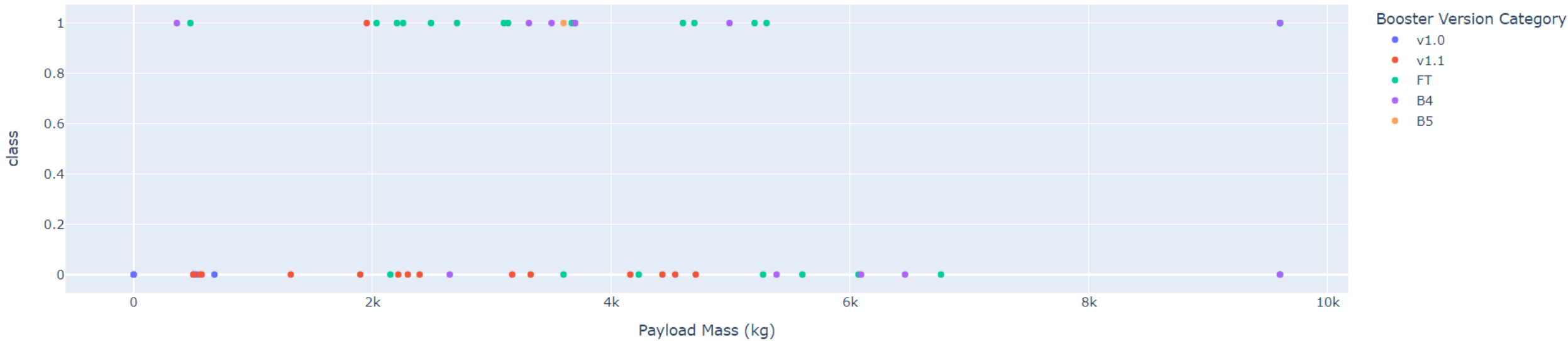


- KSC LC-39A
- CCAFS LC-40
- VAFB SLC-4E
- CCAFS SLC-40

Payload range (Kg):



All sites - payload mass between 0kg and 9,600kg



Section 5

Predictive Analysis (Classification)

Classification Accuracy

```
algorithms = {'KNN':knn_cv.best_score_,
              'Tree':tree_cv.best_score_,
              'LogisticRegression':logreg_cv.best_score_}

bestalgorithm = max(algorithms, key=algorithms.get)

print('Best Algorithm is',bestalgorithm,'with a score of',algorithms[bestalgorithm])
if bestalgorithm == 'Tree':
    print('Best Params is :',tree_cv.best_params_)
if bestalgorithm == 'KNN':
    print('Best Params is :',knn_cv.best_params_)
if bestalgorithm == 'LogisticRegression':
    print('Best Params is :',logreg_cv.best_params_)
```

Best Algorithm is Tree with a score of 0.9035714285714287

Best Params is : {'criterion': 'gini', 'max_depth': 4, 'max_features': 'sqrt', 'min_samples_leaf': 2, 'min_samples_split': 5, 'splitter': 'random'}

Confusion Matrix



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!

