



IBM Developer
SKILLS NETWORK

Winning Space Race with Data Science

Essam Fathi Mohamed Sabbah
3/10/2022



Outline

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

Executive Summary

- Summary of methodologies
 - Collection of Falcon 9 launch data via SpaceX API, and web scrapping, then store it into IBM Db2 for analytics.
 - Data cleaning, replacing missing values with either mean of that table or 0.
 - Visualization of the data
 - Finally, performance of a predictive analysis to train and evaluate a best model and give prediction over the success of future landings.
- Summary of all results
 - When Payload was greater than 7500 kg falcon rocket had a higher chance of successful landing.
 - Among 11 orbit types ES L1, GEO, HEO, SSO were 100 successful with less than 6000 kg payload.
 - SpaceX has 4 launch sites, one is near California, the other three is near Florida and South Texas. All the sites are in near proximity to ocean and all the sites are bit far away from the city.
 - All models performed similarly when trained with the data at hand.

Introduction

- Project background and context
- SpaceX is a company that aims to make commercial space travel more affordable for everyone.
- This company can launch rockets for a cost of around 60 million dollars. In contrast, other providers require 165 million dollars for one launch. This is due to the fact that SpaceX can reuse the first stage of the rocket Falcon9.
- The primary cost saving agent is the high success rate of stage 1 landing and thus its reusability in future launches.
- The challenge here is to set a right costing forecast of the rocket launches through predicting its potential to land stage 1 successfully.
- Problems you want to find answers
- What are features that contribute the most to predict whether the stage one of the rocket will land successfully?
- Can we predict if new launches will be successful based on our trained model. What will be the accuracy of our predictions?

Section 1

Methodology

Methodology

Executive Summary

- **Data collection methodology:**
 - Data was collected using the SpaceX REST API and Wikipedia Web scrapping using python BeautifulSoup.
- **Perform data wrangling**
 - Data was processed using python pandas and numpy library.
- **Perform exploratory data analysis (EDA) using visualization and SQL**
- **Perform interactive visual analytics using Folium and Plotly Dash**
- **Perform predictive analysis using classification models**
 - Classification models (Experiment usability and compatibility of SVM, Tree maps, KNN, Logistic
 - Regression optimizing parameters) were built, evaluated and tuned using sklearn.

Data Collection

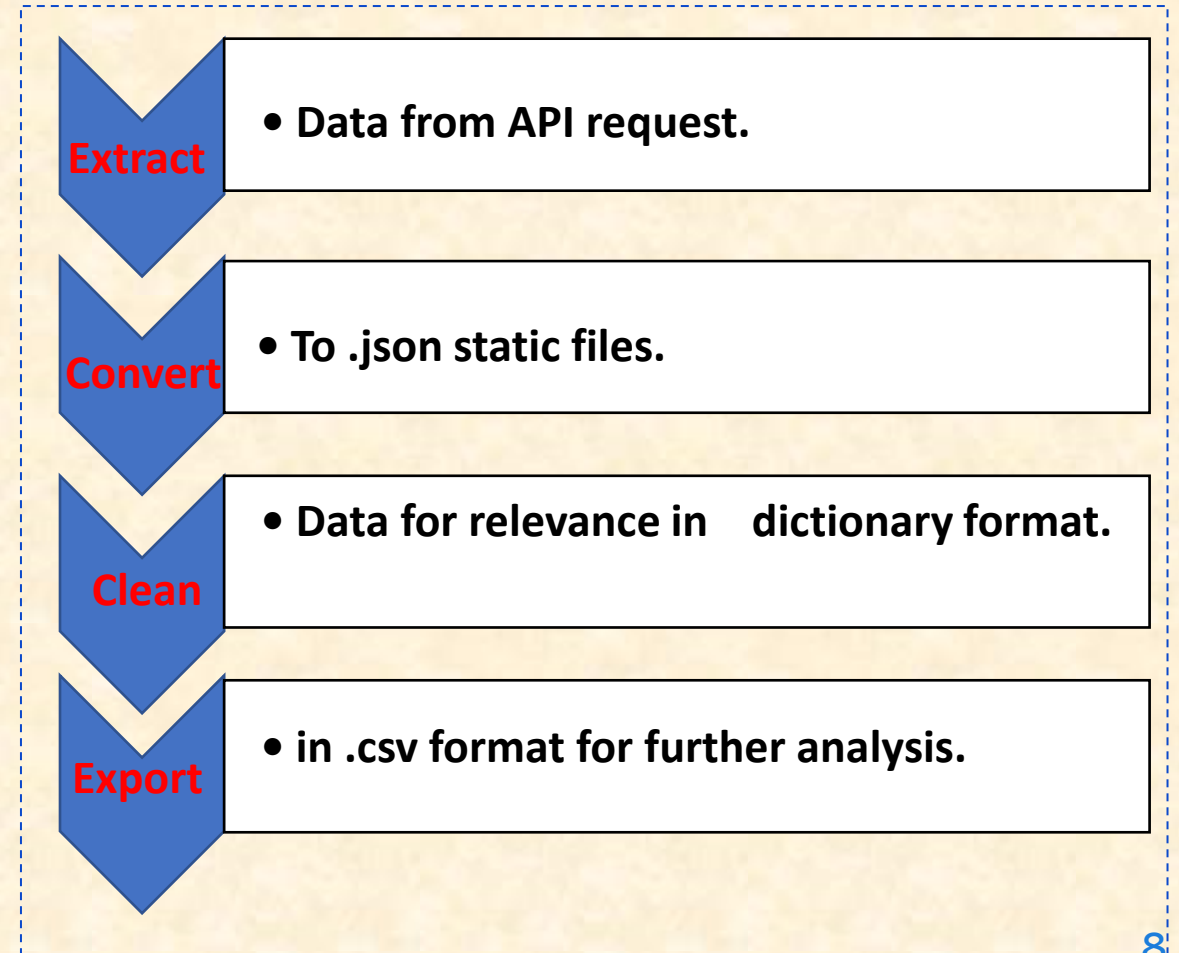
- The data is collected:

1. Using SpaceX API.
2. 2. Using BeautifulSoup library to scrap data from the Wikipedia page.

Data Collection – SpaceX API

- Request launch data from SpaceX URL using given API.
- Extract the data from the response.
- Pre-process and construct the data.
- Store the data in CSV file.

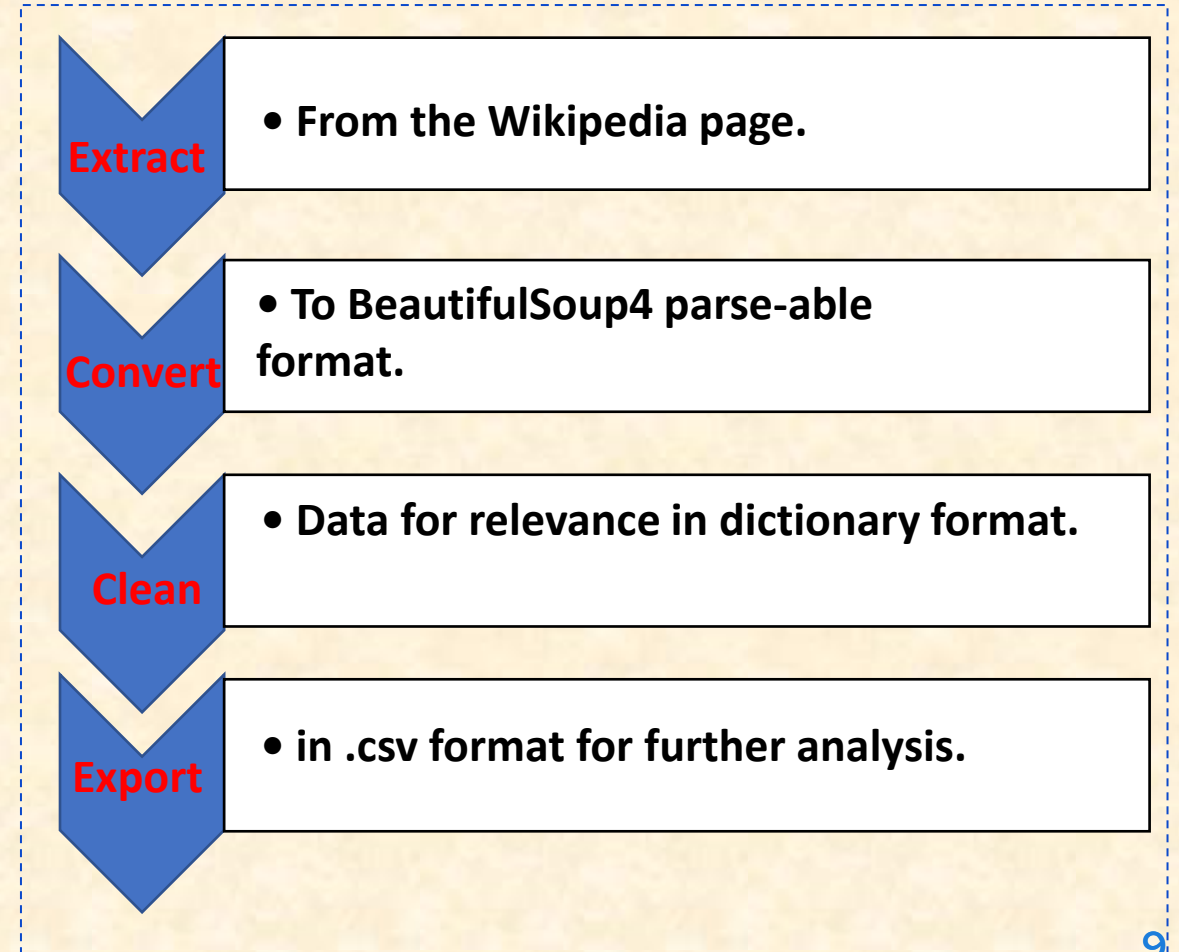
[GitHub](#)



Data Collection - Scraping

- Request for the Wikipedia page.
- Parse the table data from html text using BeautifulSoup4 library.
- Create pandas data frame from table data.
- Construct the data and store it in CSV.

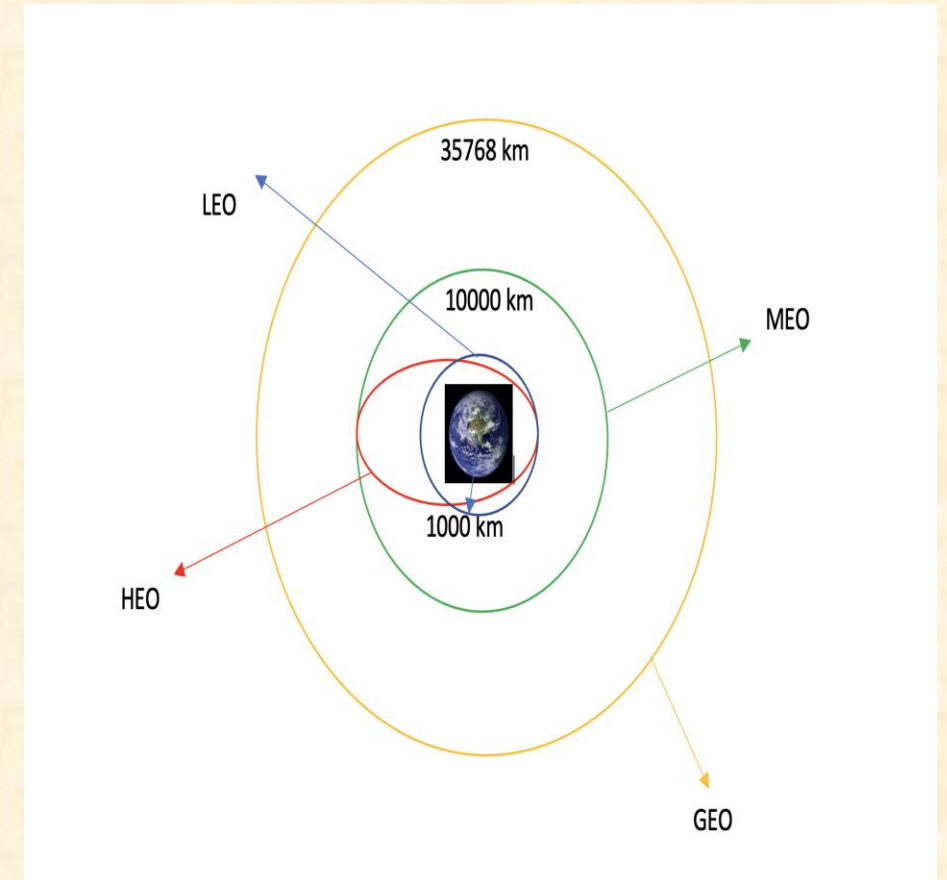
GitHub



Data Wrangling

1. Calculate the number of launches per site.
2. Number of occurrence of each orbit.
3. Number of occurrences of outcome per orbit.
4. Create landing outcome label.

[GitHub](#)



Different orbits in which flight can launched.

EDA with Data Visualization



EDA with SQL

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

Build an Interactive Map with Folium

- **Map objects** which are created and added to the folium map are given below:
 - **Markers:** Added to mark a specific area with a text label on a specific coordinate.
 - **Circles:** Added to highlight circle areas with a text label on a specific coordinate.
 - **Marker Cluster:** Marker clusters were used to simplify the containing many markers having the same coordinates.
 - **Mouse Position:** Used to get coordinate for a mouse over a point on the map (proximities). It helps to find the coordinates easily of any points of interests while exploring the map.
 - **Polyline:** It draws polyline overlays on a map. It was used to denote the distance between a launch site and its proximities(such as Railway station, city, etc.).

GitHub

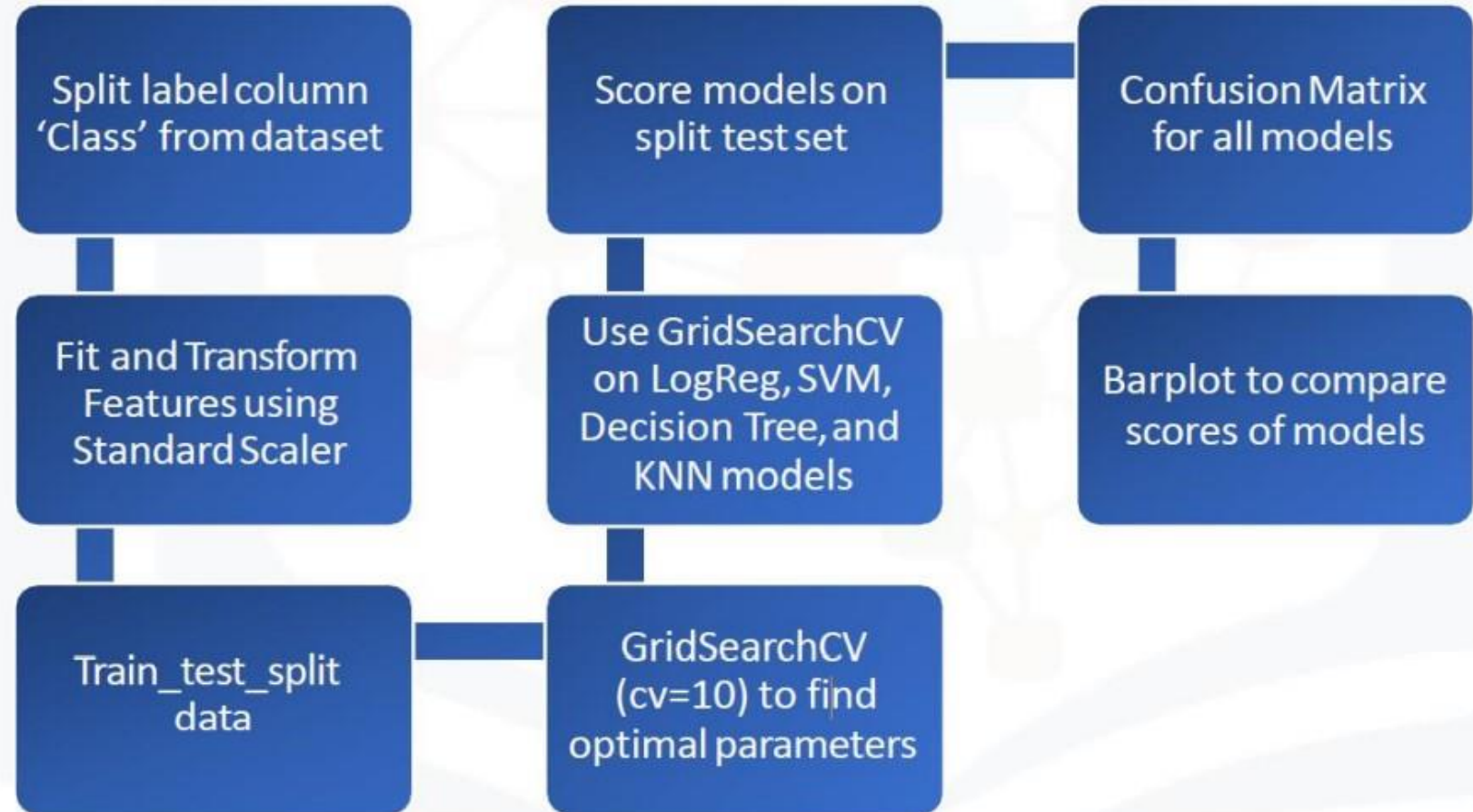
Build a Dashboard with Plotly Dash

- **Pie chart (total launches for a selected site or the total sites collection)**
 - Shows relative proportions of different sites successful landing distribution.
 - Shows % of success vs. failure for a given site.
- **Scatter Plot**
 - Showing the correlation between Outcome and Payload Mass(Kg) for different Booster Versions with freedom of selection of the range of payload mass of Interest.

[GitHub](#)

Predictive Analysis (Classification)

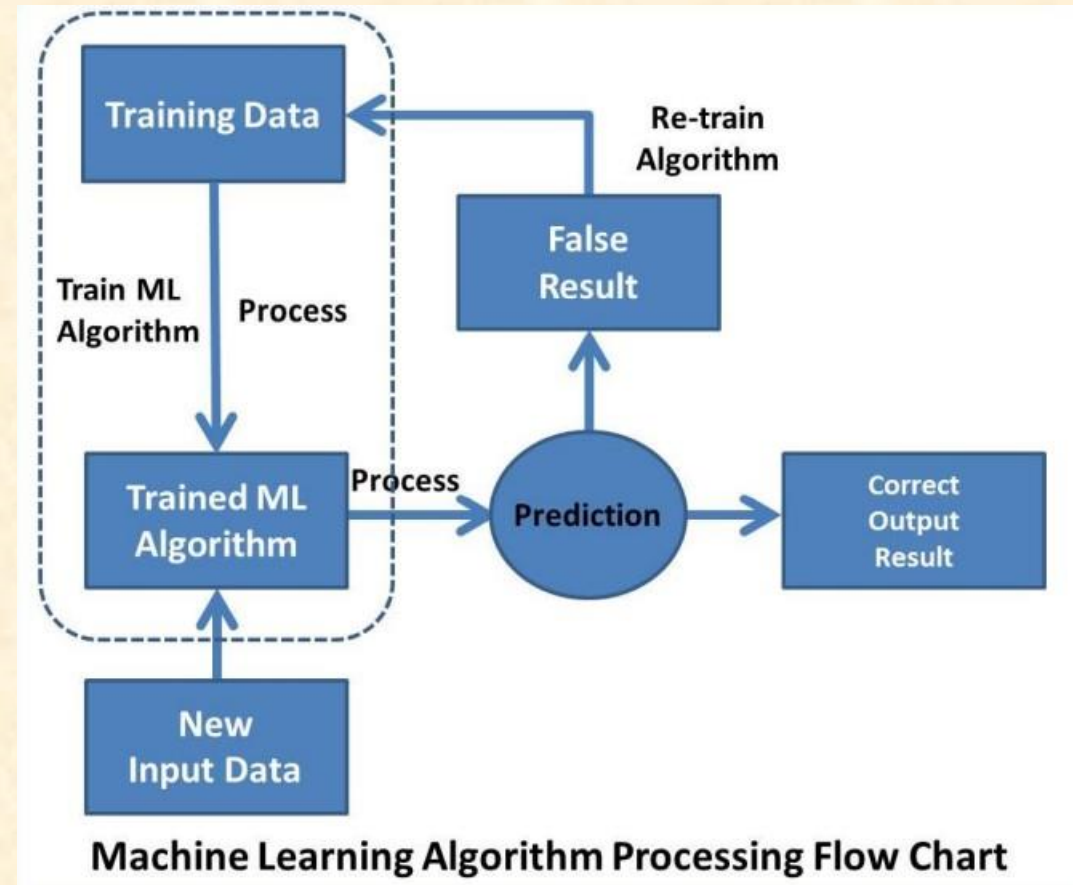
GitHub



Results

- Models were built using Scikit - Learn, data were previously normalized and models hyper parameters were found using a Grid Search with a 10 fold cross validation, in the end the best performing model has been selected based on accuracy .

[GitHub](#)

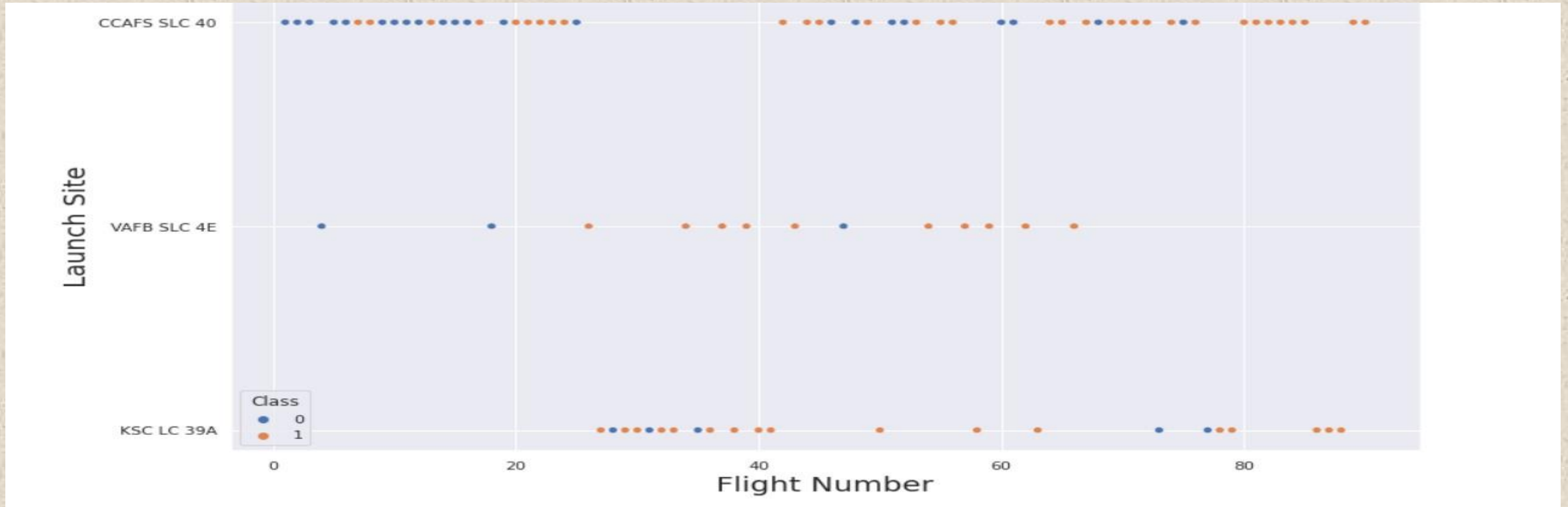


The background of the slide is an abstract composition. It features a dark blue base color. Overlaid on this are numerous diagonal streaks in shades of red and cyan. A faint, light blue grid pattern is also visible, particularly in the lower-left quadrant. The overall effect is dynamic and technological.

Section 2

Insights drawn from EDA

Flight Number vs. Launch Site



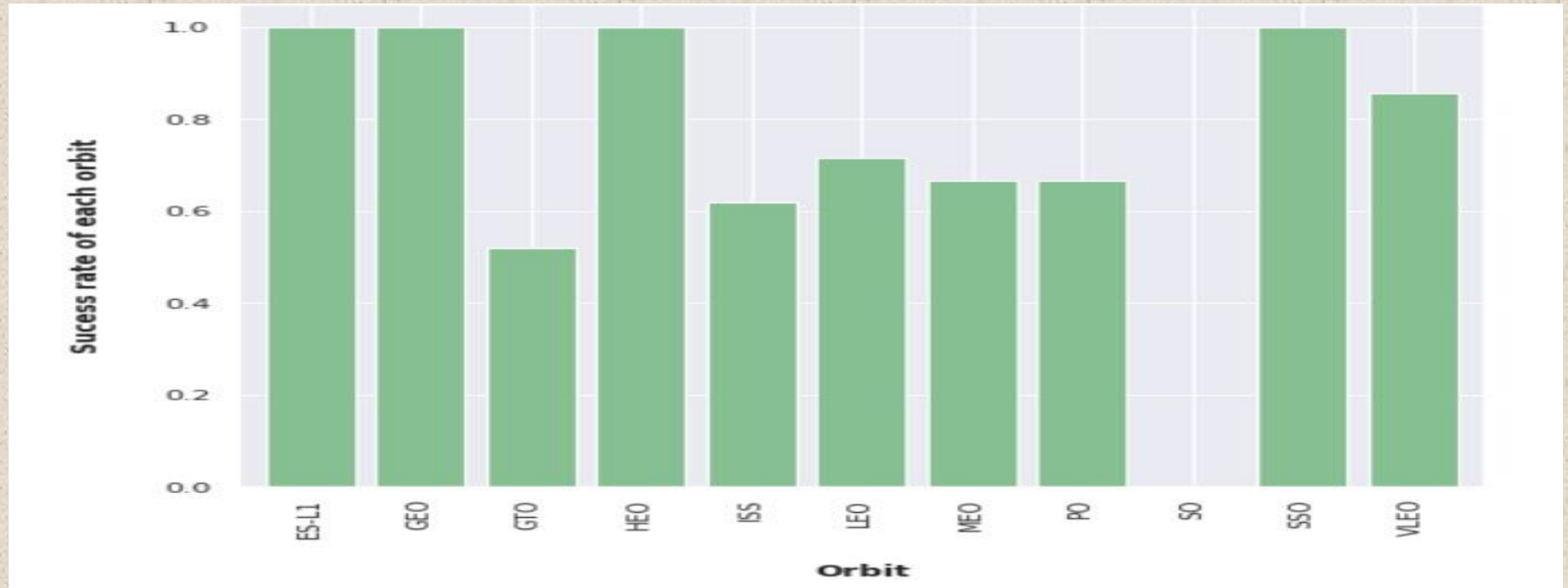
- This chart seems to indicate that a “young” launching site will probably a lower success rate than one which had a lot of rocket launched from.

Payload vs. Launch Site



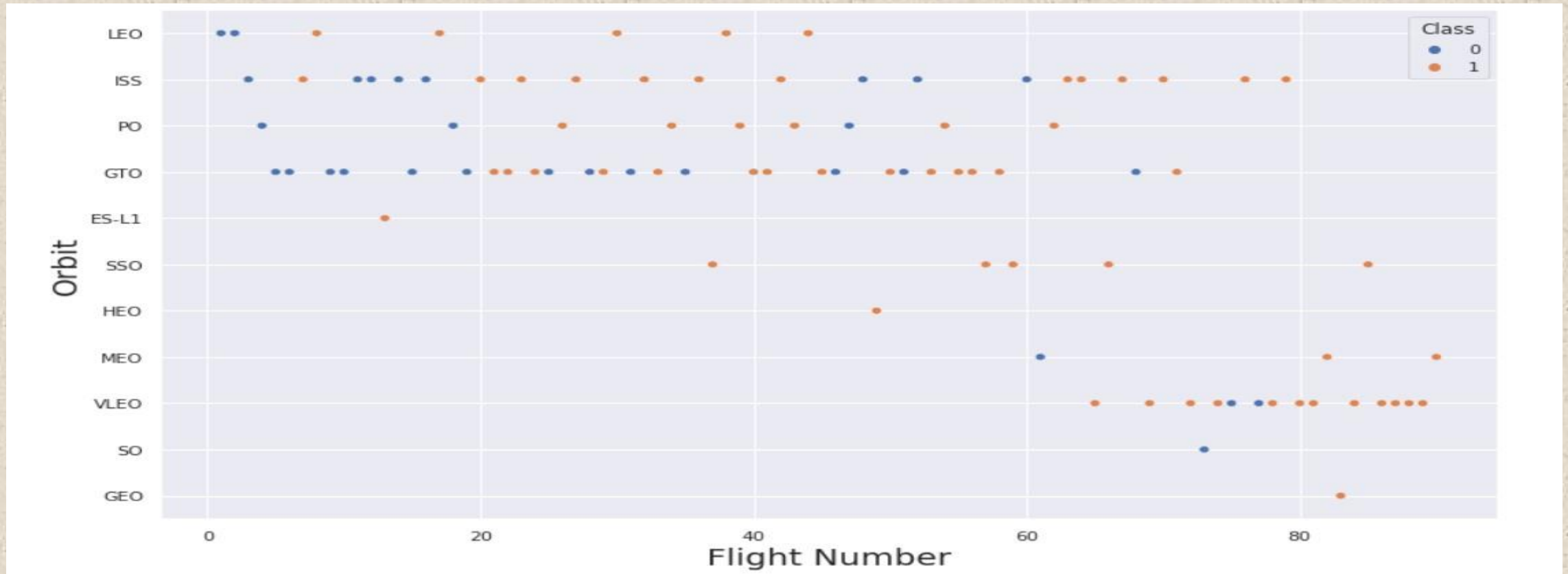
It seems like a lot of rocket launched had a payload between 500kg and 6000kg. Also, the launching site VAFB SLC 4E seems to be a site where there are not that much rocket launched. An impact of the payload could be possible but it will need further analysis.

Success Rate vs. Orbit Type



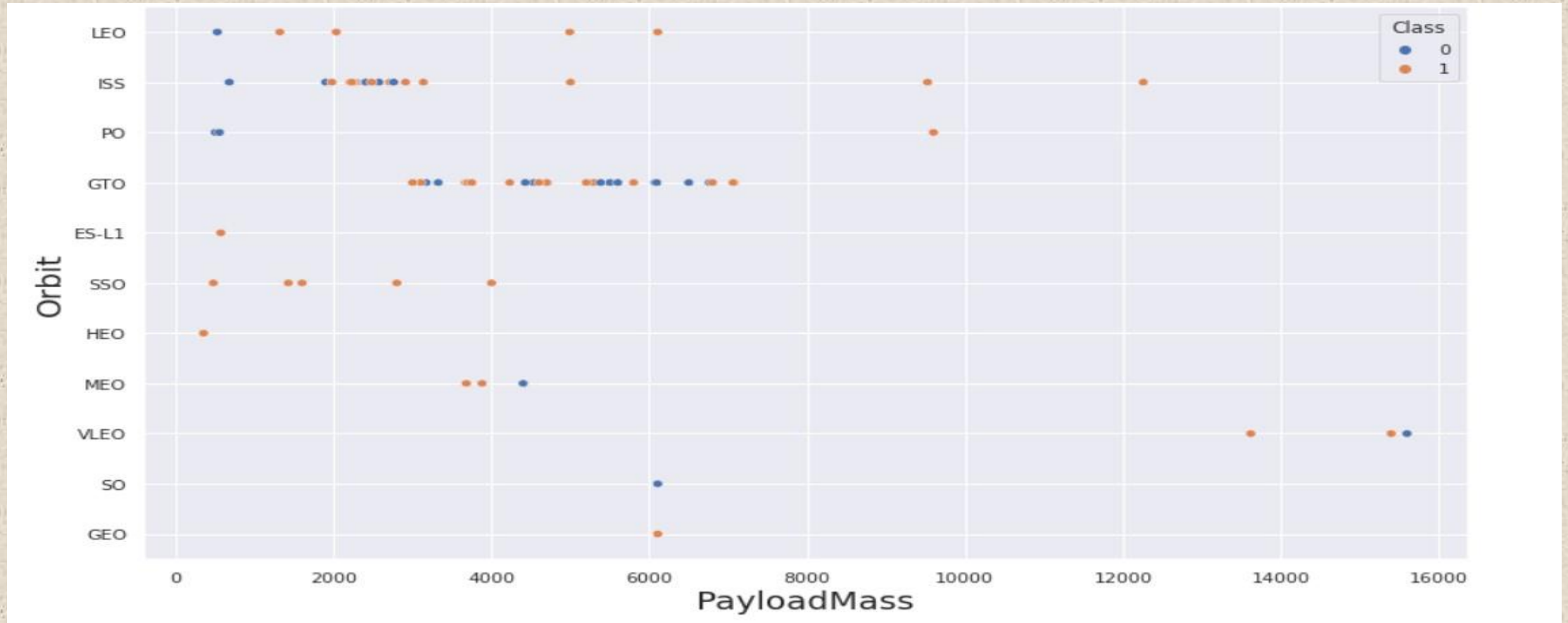
There is a strong correlation between these two indeed as we can observe the SO or GTO Orbit type are quite risky as the success rate is below 0.6. However, some orbit type provide a 1.0 success rate which is perfect but can hide suspicious data. Indeed if for this orbit type only one rocket has been launched the reliability of this hypothesis is null.

Flight Number vs. Orbit Type



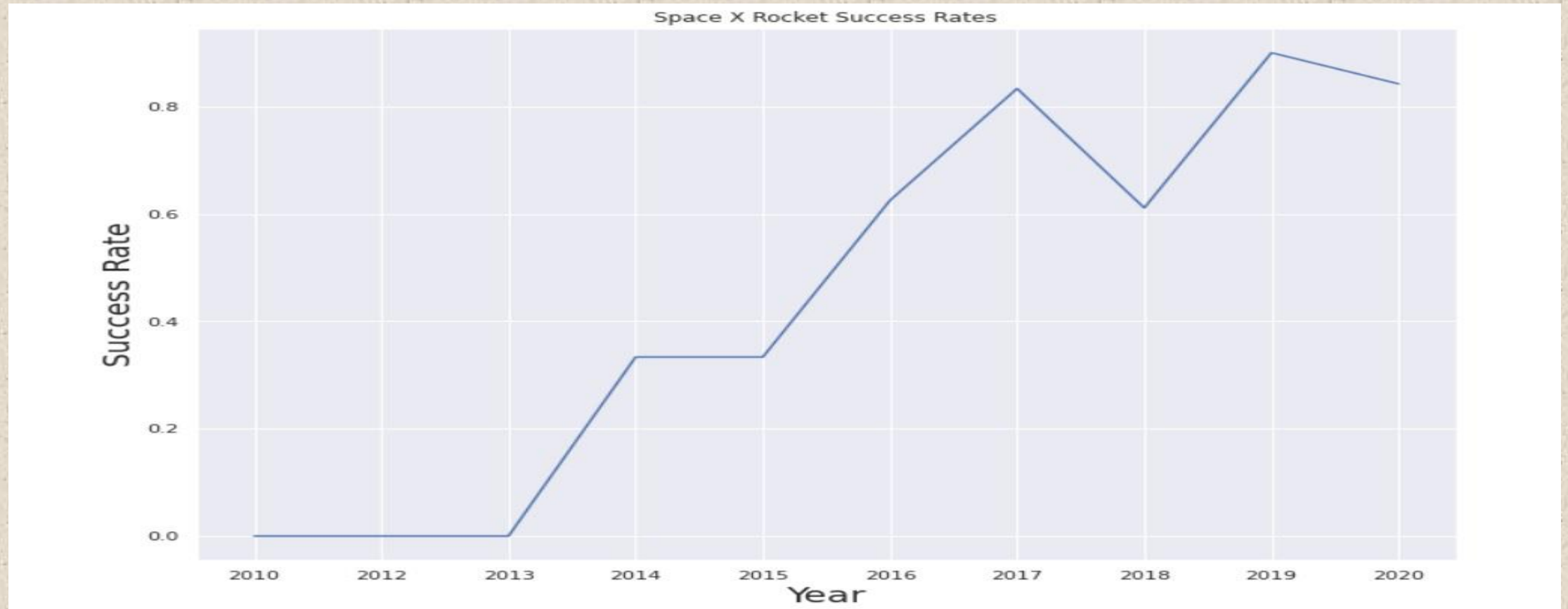
- This chart confirmed what has been said before, some Orbit type have only couples of Flights in their history and thus make data quite confusing. However for the GTO,VLEO and ISS it seems like there are enough data to be confident on those data.

Payload vs. Orbit Type



- Here we can observe that certain sites have a strong relation with the payload mass, for example the GTO and ISS.

Launch Success Yearly Trend



- Here the chart demonstrates that as Humans learn more and more through the years thanks of Sciences, it results in a significant rocket launches success rate increasing.

All Launch Site Names

- Launch sites are :
 - CCAFS LC-40
 - CCAFS SLC-40
 - KSC LC-39A
 - VAFB SLC-4E
- SQL QUERY: `select distinct(launch_site) from SPACEXTBL;`

Launch Site Names Begin with 'CCA'

- Find 5 records where launch sites begin with 'CCA'

Out[16]:	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

- SQL QUERY: `select * from SPACEXTBL where launch_site like 'CCA%' limit 5;`

Total Payload Mass

- The total payload carried by boosters from NASA is 45596 kg.
- SQL QUERY: `select sum(payload__mass__kg_) from SPACEXTBL where customer = 'NASA (CRS)';`

Average Payload Mass by F9 v1.1

- The average payload mass carried by booster version F9 v1.1 is 2534 kg.
- SQL QUERY: `select avg(payload_mass__kg_) from SPACEXTBL where booster_version like 'F9 v1.1%';`

First Successful Ground Landing Date

- The dates of the first successful landing outcome on ground pad was 2015-12-22.
- SQL QUERY: `select min(DATE) from SPACEXTBL where landing__outcome = 'Success (ground pad)';`

Successful Drone Ship Landing with Payload between 4000 and 6000

- List the names of boosters which have successfully landed on drone ship and had payload mass greater than 4000 but less than 6000:
 - F9 FT B1032.1
 - F9 B4 B1040.1
 - F9 B4 B1043.1
 - SQL QUERY: `select distinct(booster_version) from SPACEXTBL where landing__outcome = 'Success (drone ship)' AND payload_mass__kg_ > 4000 AND payload_mass__kg_ < 6000;`

Total Number of Successful and Failure Mission Outcomes

- Calculate the total number of successful and failure mission outcomes:

```
Out[23]:
```

mission_outcome	total_number
Failure (in flight)	1
Success	99
Success (payload status unclear)	1

- SQL QUERY: `select mission_outcome,count(*) from SPACEXTBL group by mission_outcome;`

Boosters Carried Maximum Payload

- List the names of the booster which have carried the maximum payload mass

Out[24]: **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
F9 B5 B1051.4
F9 B5 B1051.6
F9 B5 B1056.4
F9 B5 B1058.3
F9 B5 B1060.2
F9 B5 B1060.3
```

- SQL QUERY: `select distinct(booster_version) from SPACEXTBL where payload_mass__kg_ = (select max(payload_mass__kg_) from SPACEXTBL);`

2015 Launch Records

- List the failed landing_outcomes in drone ship, their booster versions, and launch site names for in year 2015

```
Out[25]: landing_outcome booster_version launch_site
Failure (drone ship)    F9 v1.1 B1012 CCAFS LC-40
Failure (drone ship)    F9 v1.1 B1015 CCAFS LC-40
```

- **SQL QUERY:** `select landing__outcome,booster_version,launch_site,date from SPACEXTBL where Date like '2015%' and landing__outcome = 'Failure (drone ship)';`

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

- Rank 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

Out[26]:

landing__outcome	total_number
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

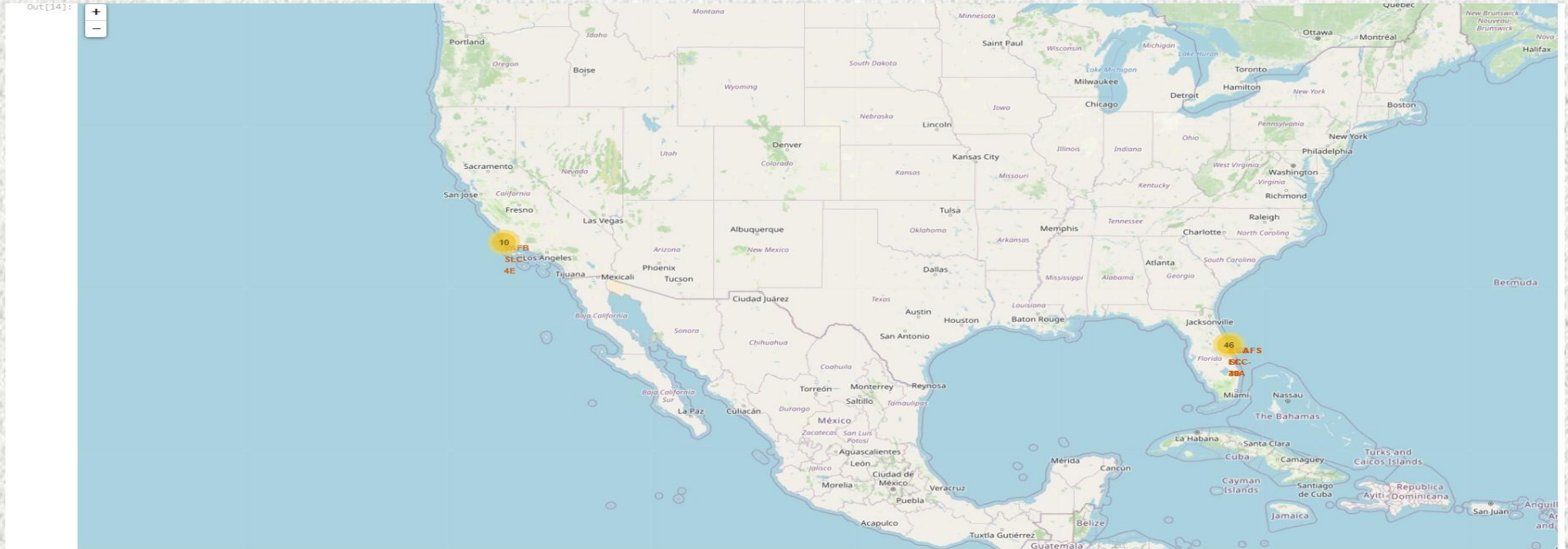
- SQL QUERY:** `select landing__outcome,count(landing__outcome) as count from SPACEXTBL where DATE BETWEEN '2010-06-04' and '2017-03-20' GROUP BY landing__outcome ORDER BY count(landing__outcome) DESC;`

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

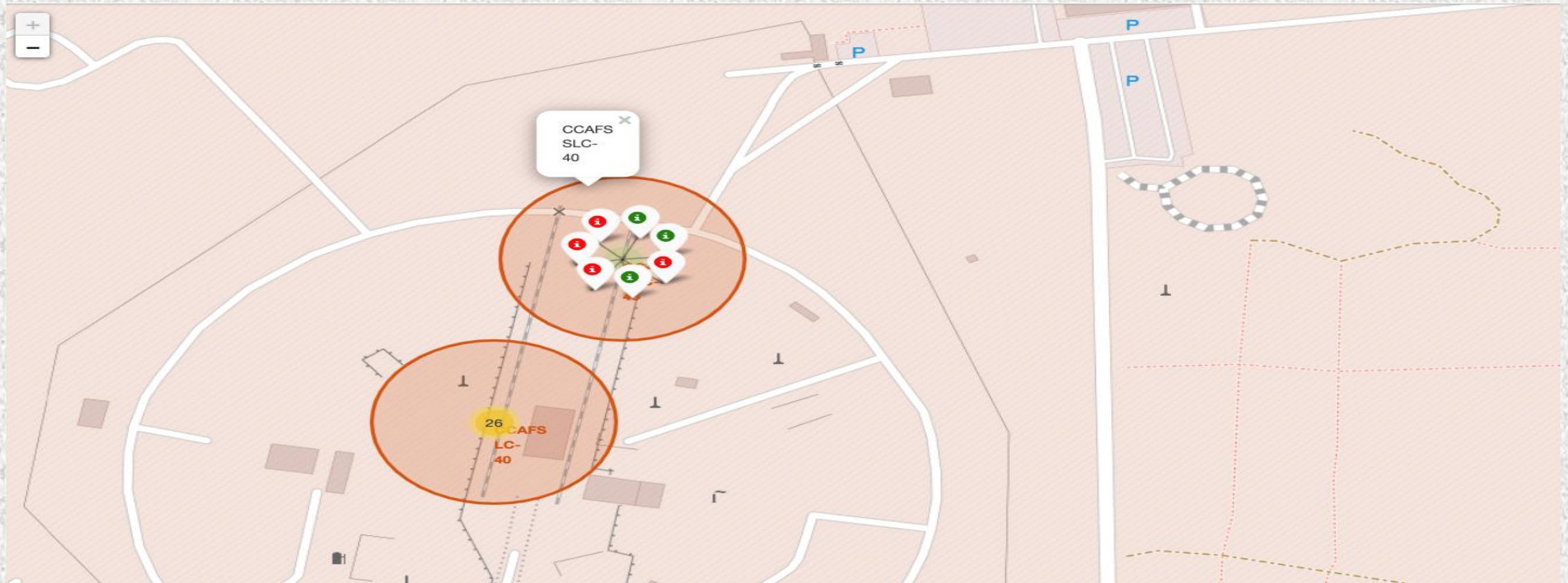
Launch Sites Proximities Analysis

SPACEX LAUNCH SITES



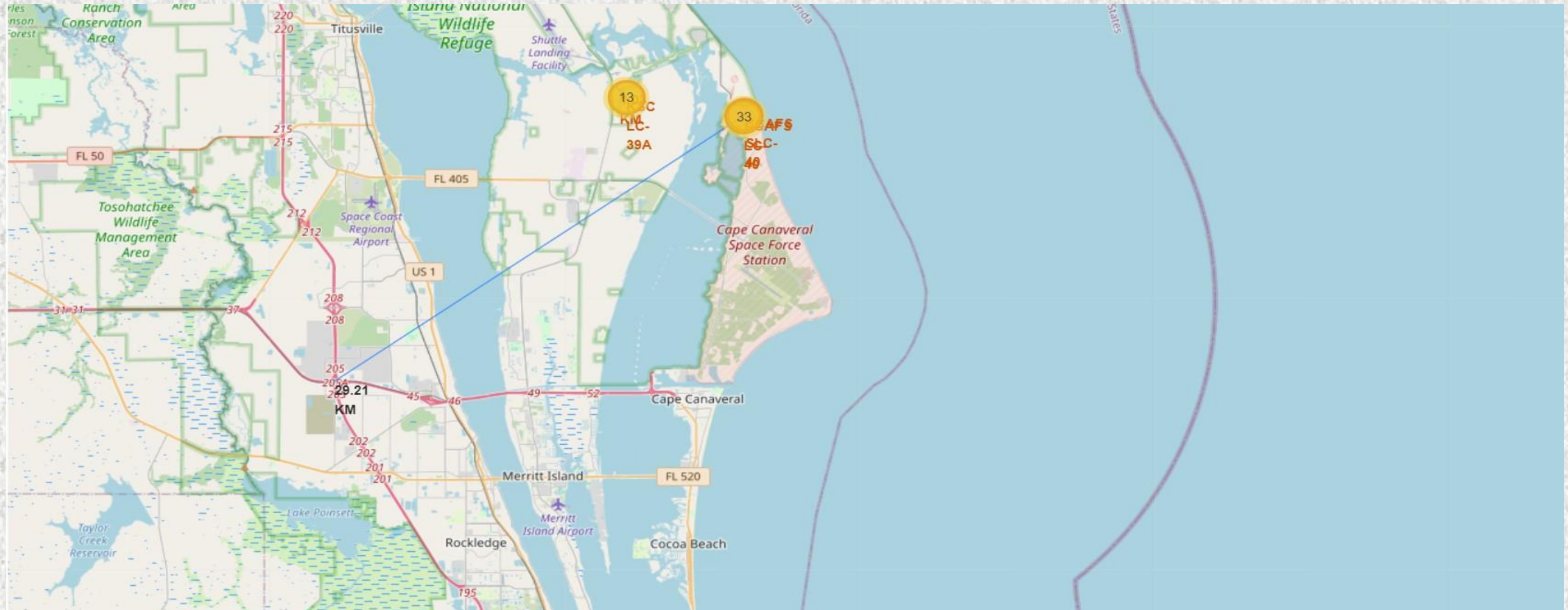
- Here we can observe launching sites in the US marked in orange, it is a little bit hard to see in Florida as the three sites are very close.

SUCCESSFUL AND FAILED LAUNCHES



- Marker clusters is used to simplify the map containing many markers having the same coordinate.
- Successful launches are marked using a green marker and failed launches are marked using a red marker.

DISTANCE FROM LAUNCH SITE TO NEAREST RAILWAY STATION



- The screenshot show the distance between launch site and nearest railway station.



Section 4

Build a Dashboard with Plotly Dash

SUCCESSFUL LAUNCHES BY SITE

Success Launches for ALL SITES



- “KSC LC 39A” has the highest success ratio where as “VAFB SLC 4E” had the lowest success ratio among all four launch sites.

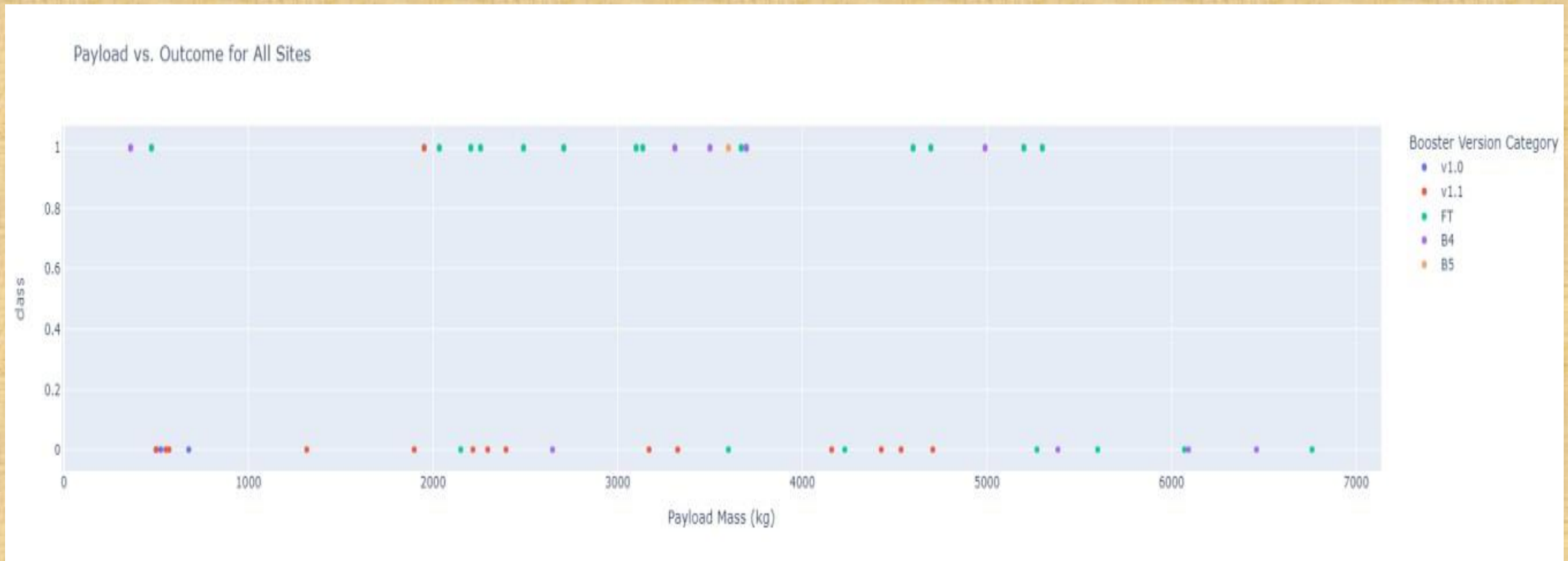
SITE WITH HIGHEST LAUNCH SUCCESS RATIO

Success Launches for site KSC LC-39A



- The KSC LC-39A has almost a 77% of success ratio and a 23% failure ratio.

PAYLOAD VS. LAUNCH OUTCOME

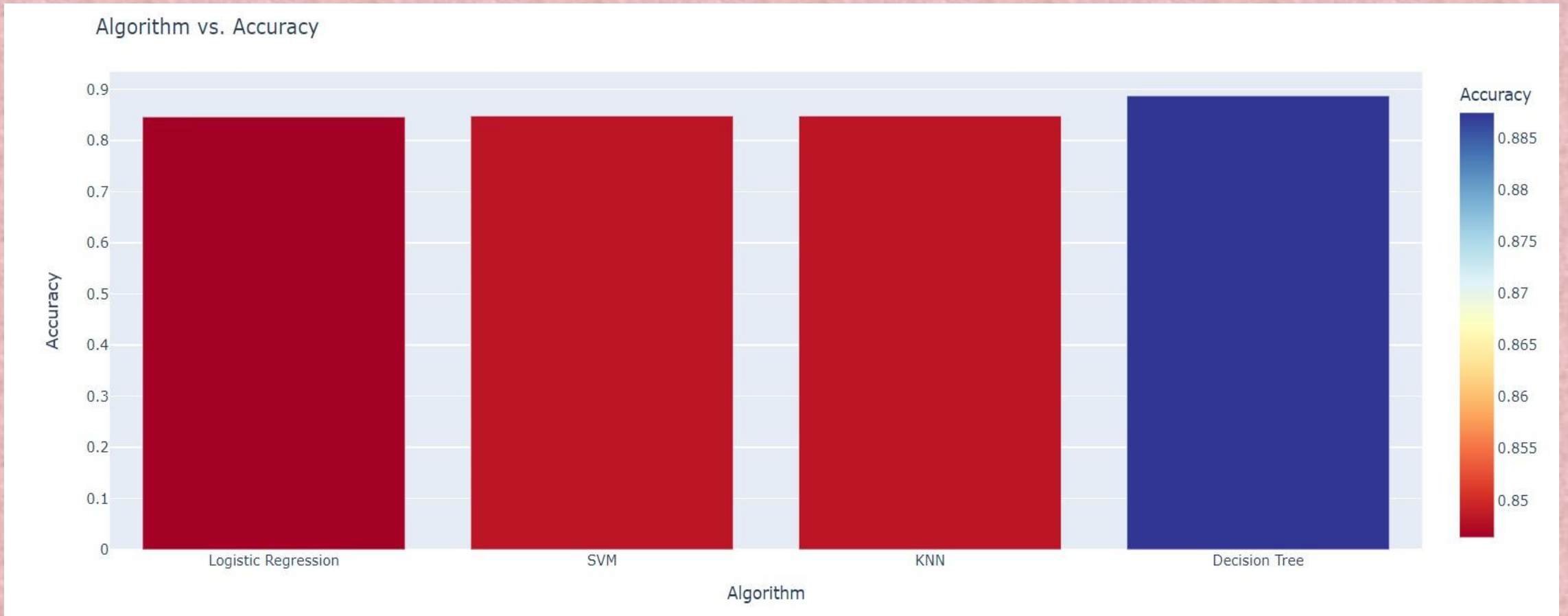


- FT booster have highest success rate while v1.1 have least success rate.

Section 5

Predictive Analysis (Classification)

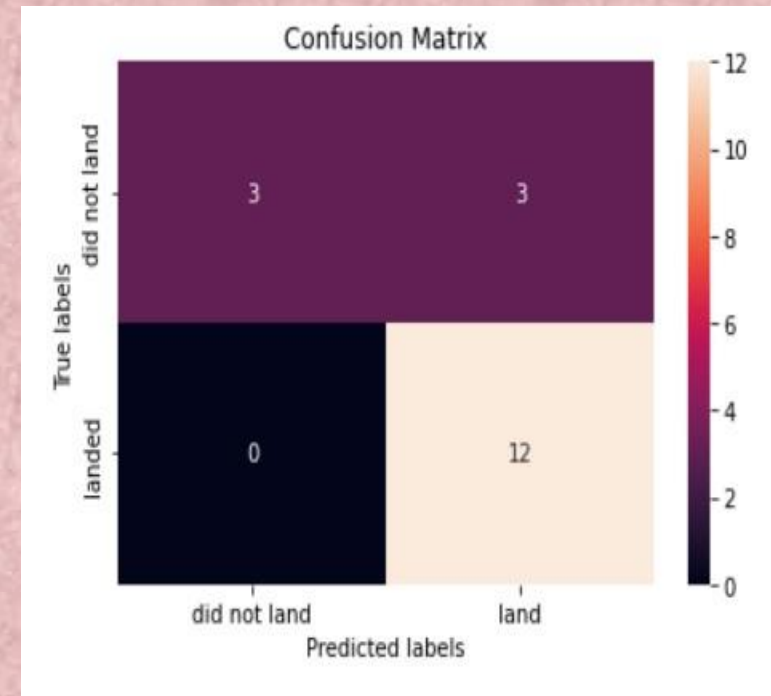
Classification Accuracy



- The best model based on the accuracy is a Decision Tree Classifier with a score of 0.884.

Confusion Matrix

- The Confusion Matrix of the Decision Tree Classifier
 - True Positive : 12
 - False Negative : 0
 - True Negative : 3
 - False Positive : 3
- The model is quite interesting as it predicts a lot of times the good labels, however 3 times it predicted the success of the mission and the mission failed. Reducing the amount of False Positive would be a good idea to avoid spending Millions and years of work.
- It could be done using Boosting or maybe look at a model with a lower accuracy but a better precision.



Conclusions

- There are many parameters when considering launching rockets in space.
- The Booster version is definitely one of this essential parameter.
- The Orbit, Payload Mass are also important.
- Machine Learning models can really helps to understand if a mission will be a success or a failure as it will learn from data of all previous launches. As we saw a model is able to predict with a high accuracy the reliability of a mission.
- However, more data would be useful to have better, I have no doubt that engineer and scientists use these data in their predictions.

Appendix

- Haversine formula
- ADGGoogleMaps Module (not used but created)
- Module sqlserver (ADGSQLSERVER)
- PythonAnywhere 24/7 dashboard

Thank you!

