



IBM Developer  
SKILLS NETWORK

# Winning Space Race with Data Science

Christian Perez  
08/30/2024



# Outline

---

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

# Executive Summary

---

- Summary of methodologies
- Summary of all results

# Introduction

---

- Project background and context
- Problems you want to find answers





## Section 1

# Methodology

# Methodology

---

## Executive Summary

- Data collection methodology:
  - Data was retrieved from SpaceX website using API and requesting data with Request library
- Perform data wrangling
  - Using Pandas and Numpy libraries to analyze and identify data like numerical and categorical and perform changes like standardization of the data.
- Perform exploratory data analysis (EDA) using visualization and SQL
- Perform interactive visual analytics using Folium and Plotly Dash
- Perform predictive analysis using classification models
  - It is needed to evaluate what alternatives of classification are in the dataset to fit what model will be used, convert categories, it is needed to calculate the distance of the points regarding to the predicted values. Evaluation of models apply to ensure the

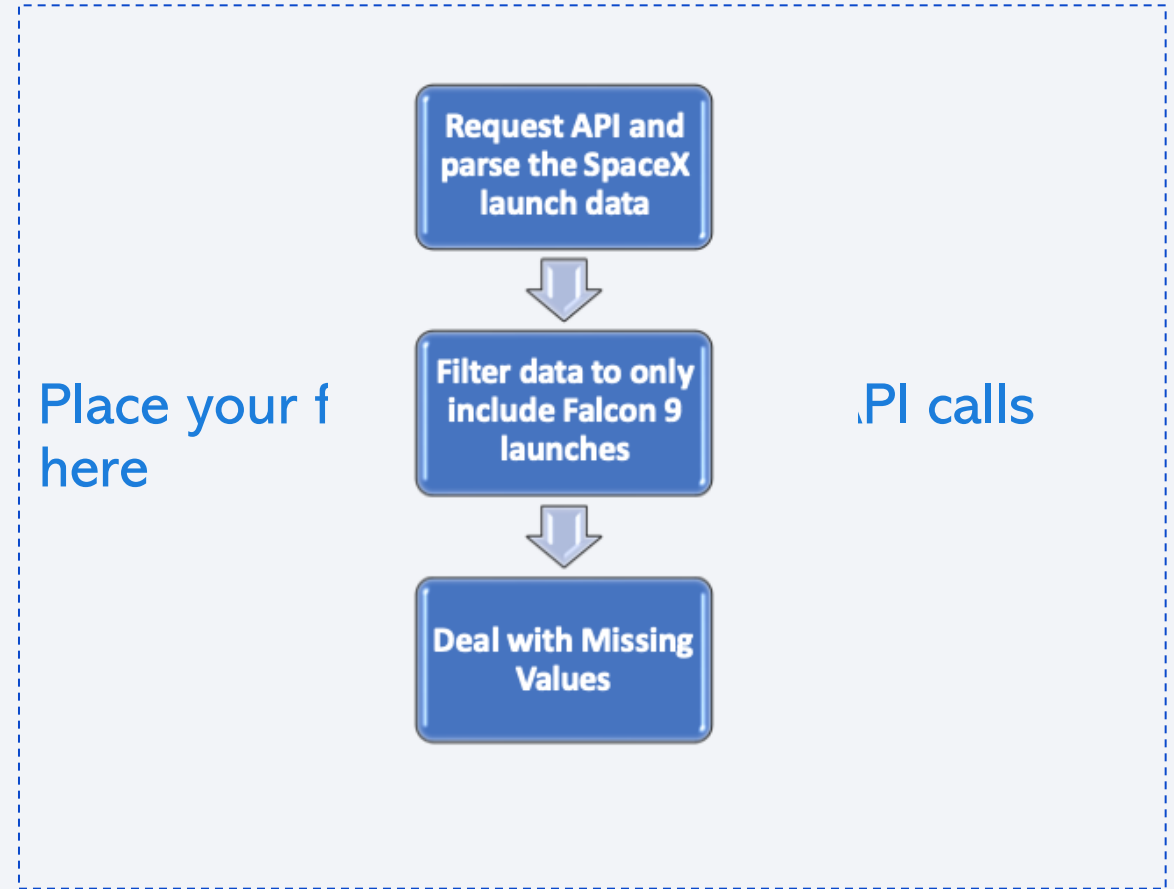
# Data Collection

---

- Describe how data sets were collected. Data was retrieved from SpaceX website using API and requesting data with Request library
- You need to present your data collection process use key phrases and flowcharts

# Data Collection – SpaceX API

- SpaceX offers a public API from where data can be obtained and then used
- <https://github.com/ChristianPerezRamirez/IBMCapstone/blob/main/jupyter-labs-spacex-data-collection-api.ipynb>





# Data Collection - Scraping

---

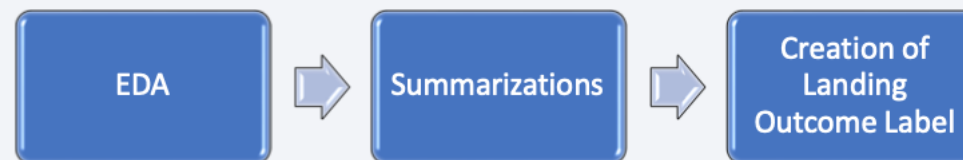
- SpaceX offers a public API from where data can be obtained and then used
- <https://github.com/ChristianPerezRamirez/IBMCapstone/blob/main/jupyter-labs-webscraping.ipynb>



# Data Wrangling

---

- Exploratory Data Analysis was performed on the dataset.
- Summary launches per site, occurrences of each orbit and occurrences of mission outcome per orbit type were calculated and the landing outcome label was created from Outcome column.



# EDA with Data Visualization

---

- Summarize what charts were plotted and why you used those charts
  - Top 5 launch sites whose name begins with the string 'CCA';
  - Total payload mass carried by boosters launched by NASA (CRS);
  - Average payload mass carried by booster version F9 v1.1;
  - Date when the first successful landing outcome in ground pad was achieved;
  - Names of the boosters which have success in drone ship and have payload mass between 4000 and 6000 kg;
  - Total number of successful and failure mission outcomes;
  - Names of the booster versions which have carried the maximum payload mass;
  - Failed landing outcomes in droneship, their booster versions, and launch site names for in year 2015; and
  - Rank of the count of landing outcomes (such as Failure (droneship) or Success (ground pad)) between the date 2010-06-04 and 2017-03-20.
- <https://github.com/ChristianPerezRamirez/IBMCapstone/blob/main/edadataviz.ipynb>

# EDA with SQL

---

- Scatterplots and bar plots were used to visualize the relationship between pair of features:
- Payload Mass X Flight Number, Launch Site X Flight Number, Launch Site X Payload Mass, Orbit and Flight Number, Payload and Orbit
- [https://github.com/ChristianPerezRamirez/IBMCapstone/blob/main/jupyter-labs-eda-sql-coursera\\_sqlite.ipynb](https://github.com/ChristianPerezRamirez/IBMCapstone/blob/main/jupyter-labs-eda-sql-coursera_sqlite.ipynb)

# Build an Interactive Map with Folium

---

- Markers, circles, lines and marker clusters were used with Folium Maps
- Markers indicate points like launch sites;
- Circles indicate highlighted areas around specific coordinates, like NASA Johnson Space Center;
- Marker clusters indicates groups of events in each coordinate, like launches in a launch site; and
- Lines are used to indicate distances between two coordinates.
- [https://github.com/ChristianPerezRamirez/IBMCapstone/blob/main/lab\\_jupyter\\_launch\\_site\\_location.ipynb](https://github.com/ChristianPerezRamirez/IBMCapstone/blob/main/lab_jupyter_launch_site_location.ipynb)

# Build a Dashboard with Plotly Dash

---

- Interactive dashboard with Plotly dash
- Pie charts showing the total launches by a certain sites
- Scatter graph showing the relationship with Outcome and Payload Mass (Kg) for the different booster version.
- <https://github.com/ChristianPerezRamirez/IBMCapstone/tree/main>



# Predictive Analysis (Classification)

---

- Data using numpy and pandas, transformed the data, split our data into training and testing.
- Different machine learning models and tune different hyperparameters using GridSearchCV.
- Accuracy as the metric for our model, improved the model using feature engineering and algorithm tuning.
- Best performing classification model.
- <https://github.com/ChristianPerezRamirez/IBMCapstone/blob/main/Machine%20Learning%20Prediction.ipynb>

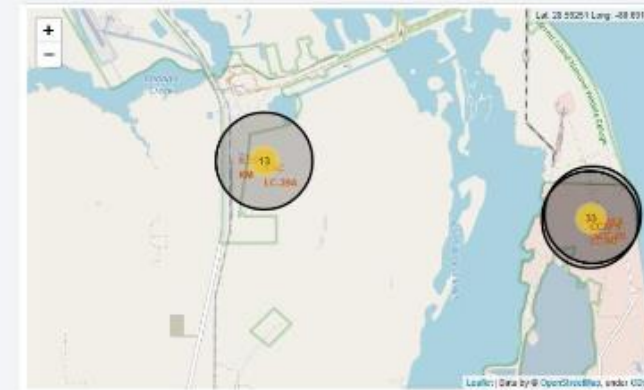
# Results

---

- Space X uses 4 different launch sites;
- First launches were done to Space X itself and NASA;
- Average payload of F9 v1.1 booster is 2,928 kg;
- First success landing outcome happened in 2015 fiver year after the first launch;
- Falcon 9 booster versions were successful at landing in drone ships having payload above the average;
- 100% of mission outcomes were successful;
- Two booster versions failed at landing in drone ships in 2015: F9 v1.1 B1012 and F9 v1.1 B1015;
- The number of landing outcomes became as better through the years.

# Results

---





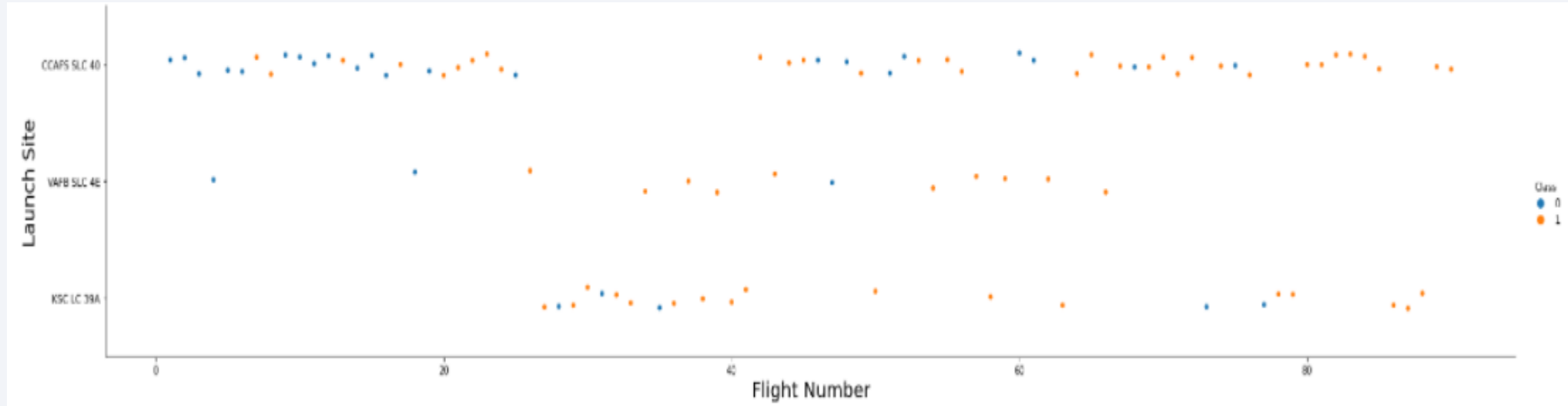


Section 2

# Insights drawn from EDA

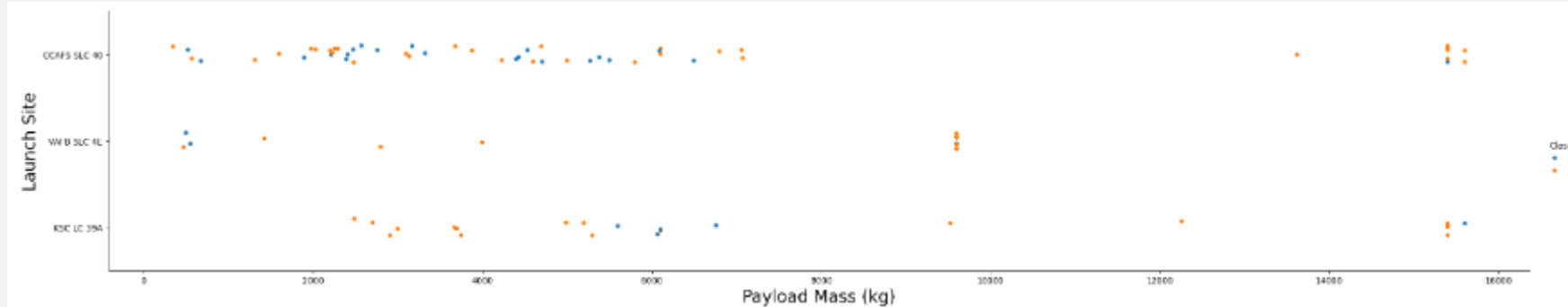


# Flight Number vs. Launch Site



The larger the flight amount at a launch site, the greater the success rate at a launch site.

# Payload vs. Launch Site



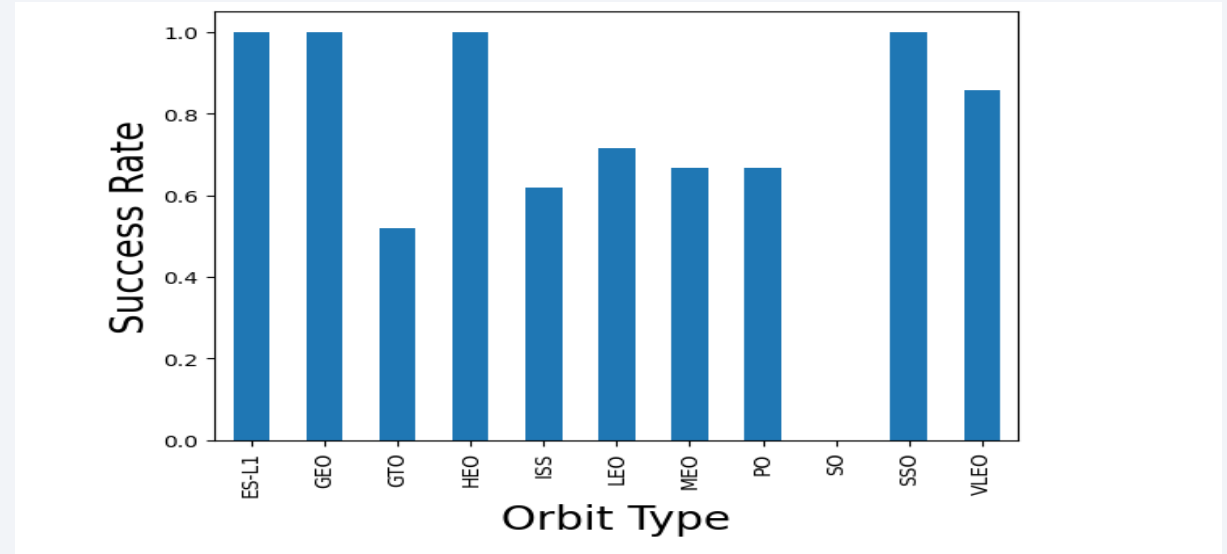
Now if you observe Payload Mass Vs. Launch Site scatter point chart you will find for the VAFB-SLC launchsite there are no rockets launched for heavypayload mass(greater than 10000).



# Success Rate vs. Orbit Type

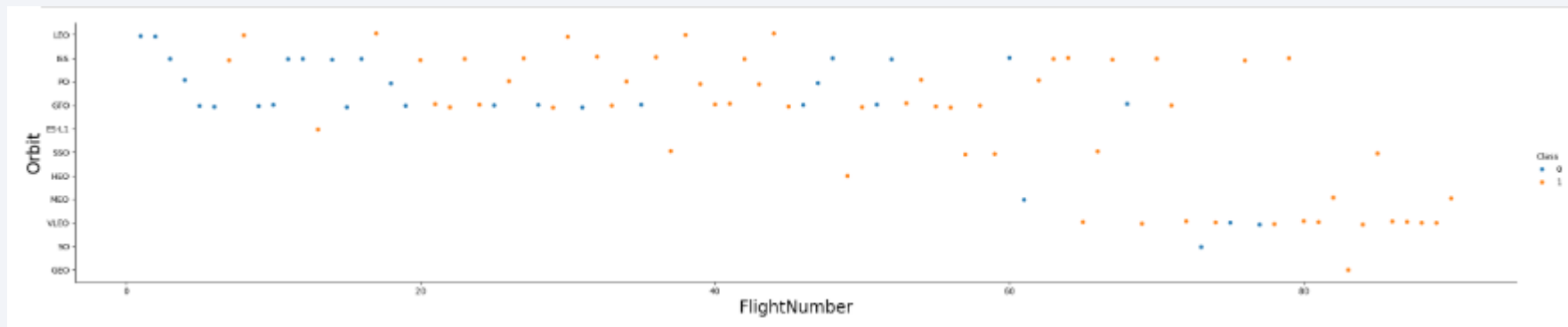
---

- ES-L1, GEO, HEO, SSO, VLEO had the most success rate.

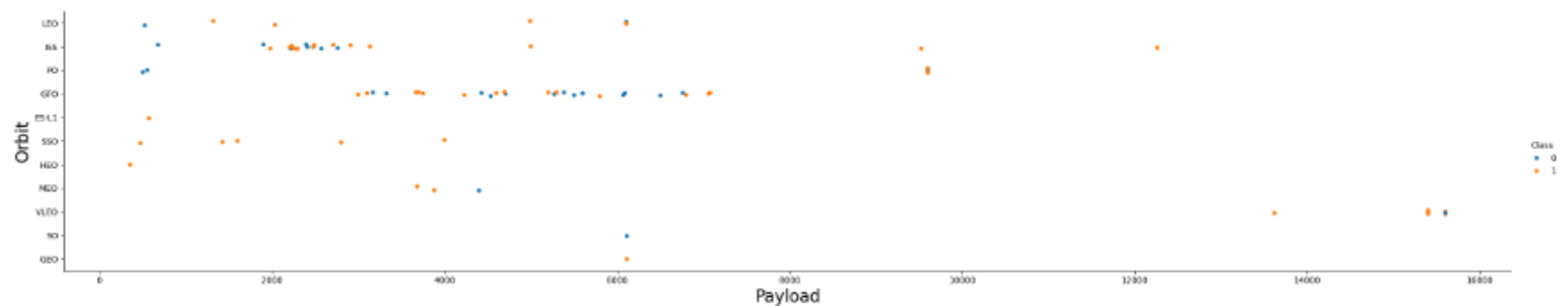


# Flight Number vs. Orbit Type

- LEO orbit, success is related to the number of flights whereas in the GTO orbit, there is no relationship between flight number and the orbit. Show the screenshot of the scatter plot with explanations



# Payload vs. Orbit Type

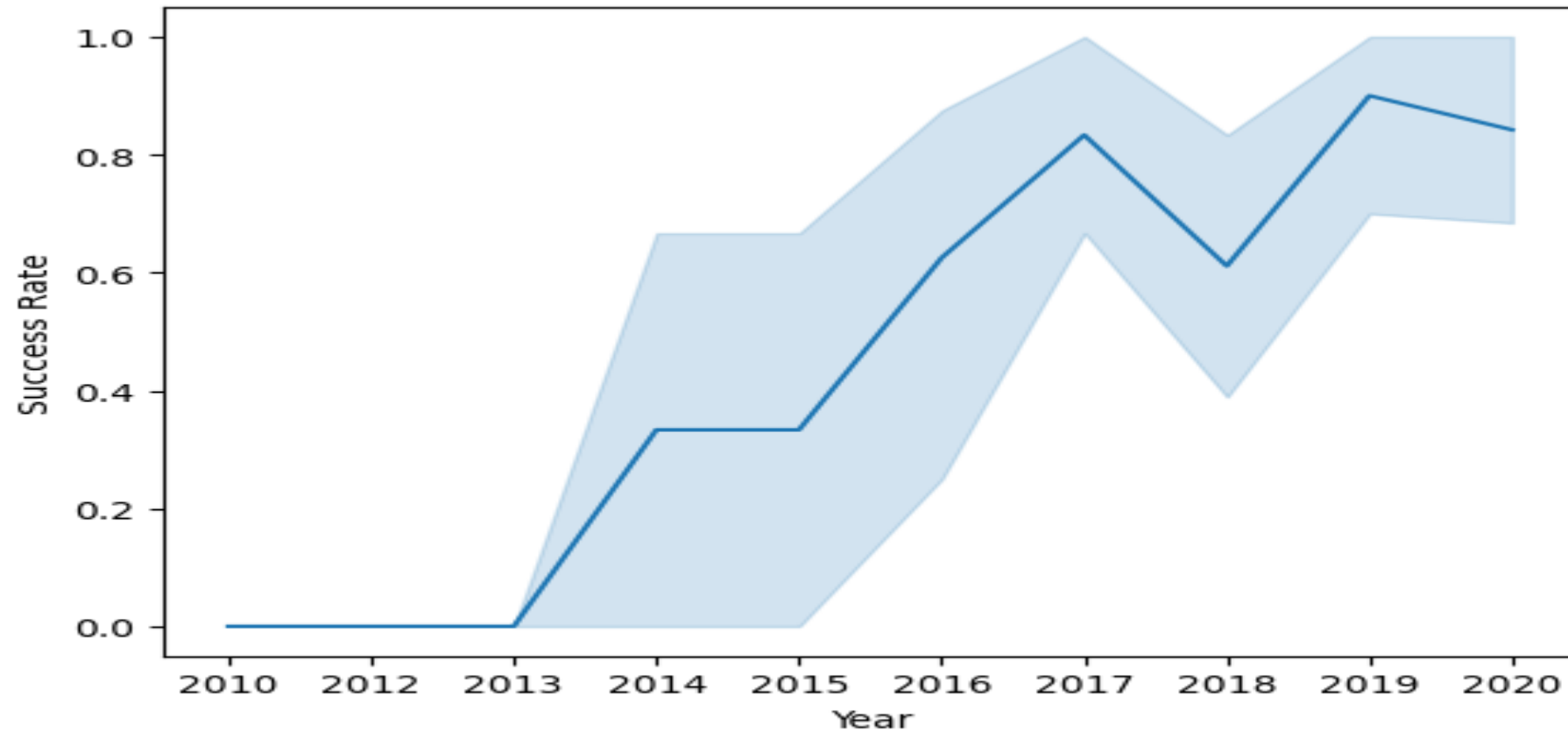


With heavy payloads the successful landing or positive landing rate are more for Polar, LEO and ISS.

However, for GEO, it's difficult to distinguish between successful and unsuccessful landings as both outcomes are present.

# Launch Success Yearly Trend

---



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

# All Launch Site Names

---

Display the names of the unique launch sites in the space mission

```
In [13]: %sql Select distinct Launch_Site from SPACEXTABLE;
```

```
* sqlite:///my_data1.db  
Done.
```

```
Out[13]: Launch_Site  
_____  
CCAFS LC-40  
VAFB SLC-4E  
KSC LC-39A  
CCAFS SLC-40
```

# Launch Site Names Begin with 'CCA'

In [15]: %sql Select \* from SPACEXTABLE where Launch\_Site like '%CCA%' limit 5;

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

Out[15]:

Date	Time (UTC)	Booster_Version	Launch_Site	Payload	PAYLOAD_MASS_KG_	Orbit	Customer	Mission_Outcome	L
2010-06-04	18:45:00	F9 v1.0 B0003	CCAFS LC-40	Dragon Spacecraft Qualification Unit	0	LEO	SpaceX	Success	F
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	F
2012-05-22	7:44:00	F9 v1.0 B0005	CCAFS LC-40	Dragon demo flight C2	525	LEO (ISS)	NASA (COTS)	Success	
2012-10-08	0:35:00	F9 v1.0 B0006	CCAFS LC-40	SpaceX CRS-1	500	LEO (ISS)	NASA (CRS)	Success	
2013-03-01	15:10:00	F9 v1.0 B0007	CCAFS LC-40	SpaceX CRS-2	677	LEO (ISS)	NASA (CRS)	Success	



# Total Payload Mass

---

## Task 3

Display the total payload mass carried by boosters launched by NASA (CRS)

```
In [16]: %sql Select sum(PAYLOAD_MASS__KG_) FROM SPACEXTABLE WHERE Customer = 'NASA (CRS)';
```

```
* sqlite:///my_data1.db
```

```
Done.
```

```
Out[16]:
```

<u>sum(PAYLOAD_MASS__KG_)</u>
45596

# Average Payload Mass by F9 v1.1

---

## Task 4

Display average payload mass carried by booster version F9 v1.1

```
In [17]: %sql Select avg(PAYLOAD_MASS__KG_) from SPACEXTABLE where Booster_Version = 'F9 v1.1';
```

```
* sqlite:///my_data1.db  
Done.
```

```
Out[17]: avg(PAYLOAD_MASS__KG_)
          2928.4
```

# First Successful Ground Landing Date

---

## Task 5

List the date when the first succesful landing outcome in ground pad was acheived.

*Hint: Use min function*

```
In [18]: %sql Select min(date) from SPACEXTABLE where Landing_Outcome = 'Success (ground pad)';
```

```
* sqlite:///my_data1.db
```

```
Done.
```

```
Out[18]: min(date)
```

```
2015-12-22
```

# Successful Drone Ship Landing with Payload between 4000 and 6000

---

## Task 6

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

In [19]: `%sql Select Payload, PAYLOAD_MASS_KG_ from SPACESTABLE where Landing_Outcome = 'Success (drone ship)' and PAYLOAD_MASS_KG_`

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

Done.

Out[19]:

Payload	PAYLOAD_MASS_KG_
JCSAT-14	4696
JCSAT-16	4600
SES-10	5300
SES-11 / EchoStar 105	5200

# Total Number of Successful and Failure Mission Outcomes

---

## Task 7

List the total number of successful and failure mission outcomes

```
In [33]: %sql SELECT MISSION_OUTCOME, COUNT(MISSION_OUTCOME) AS TOTAL_NUMBER FROM SPACEXTBL GROUP BY MISSION_OUTCOME;
```

```
* sqlite:///my_data1.db
```

```
Done.
```

```
Out[33]:
```

Mission_Outcome	TOTAL_NUMBER
Failure (in flight)	1
Success	98
Success	1
Success (payload status unclear)	1

# Boosters Carried Maximum Payload

---

## Task 8

List the names of the booster\_versions which have carried the maximum payload mass. Use a subquery

```
In [34]: %sql SELECT DISTINCT BOOSTER_VERSION FROM SPACEXTBL WHERE PAYLOAD_MASS_KG_ = (SELECT MAX(PAYLOAD_MASS_KG_) FROM SPACEXTBL)
```

```
* sqlite:///my_data1.db  
Done.
```

Out[34]: **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

---

## Task 9

List the records which will display the month names, failure landing\_outcomes in drone ship ,booster versions, launch\_site for the months in year 2015.

**Note:** SQLite does not support monthnames. So you need to use substr(Date, 6,2) as month to get the months and substr(Date,0,5)='2015' for year.

In [23]: `%sql Select substr(Date, 6,2), Landing_Outcome, Booster_Version, launch_site from SPACEXTABLE where substr(Date,0,5)='2015'`

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

Out[23]:

	substr(Date, 6,2)	Landing_Outcome	Booster_Version	Launch_Site
01	Failure (drone ship)	F9 v1.1 B1012	CCAFS LC-40	
04	Failure (drone ship)	F9 v1.1 B1015	CCAFS LC-40	

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

## Task 10

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.

```
In [40]: %sql SELECT Landing_Outcome, COUNT(Landing_Outcome) AS TOTAL_NUMBER FROM SPACEXTABLE WHERE DATE BETWEEN '2010-06-04' AND '2017-03-20'
```

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

```
Out[40]:
```

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

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 visible. The Earth's surface is dark blue, with bright yellow and orange lights from cities and towns scattered across the landmasses. The horizon line is visible, separating the dark Earth from the black space.

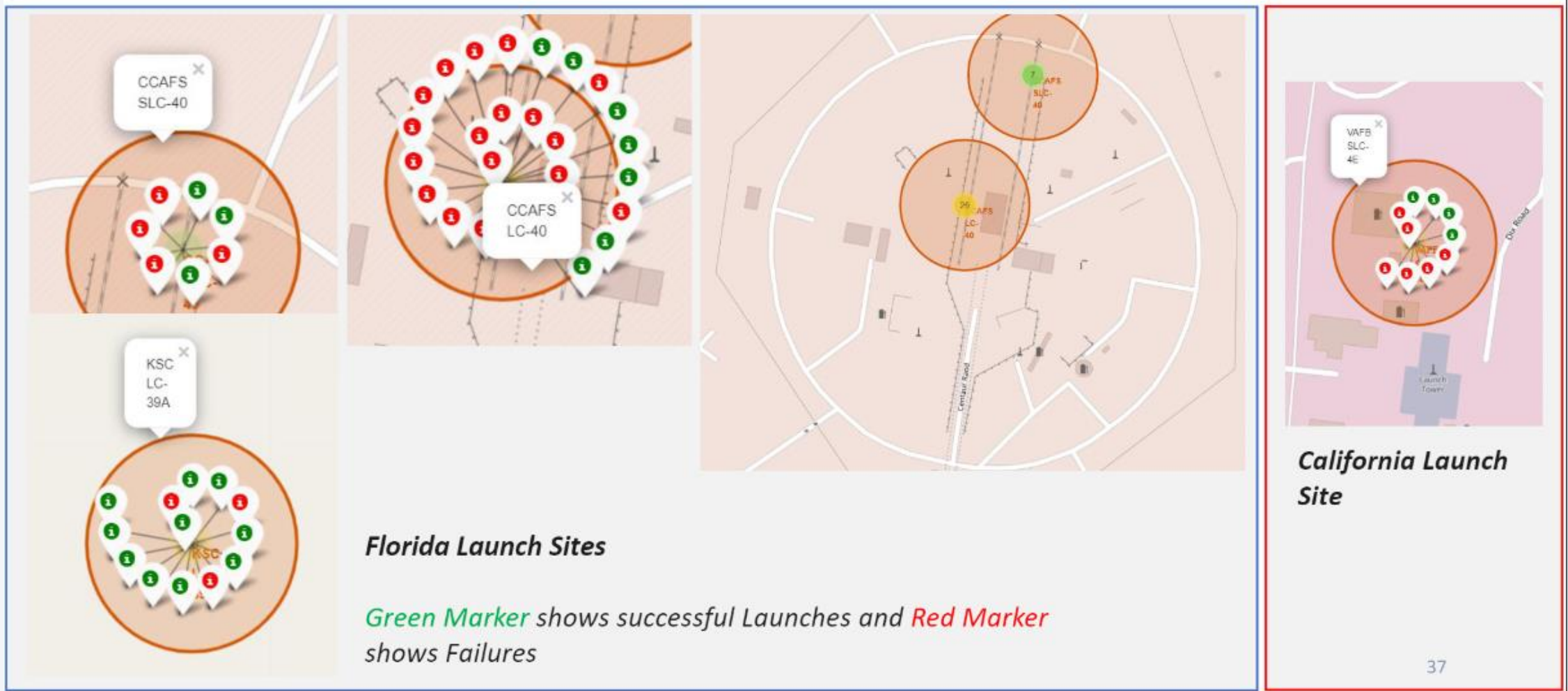
Section 3

# Launch Sites Proximities Analysis

## <Folium Map Screenshot 1>

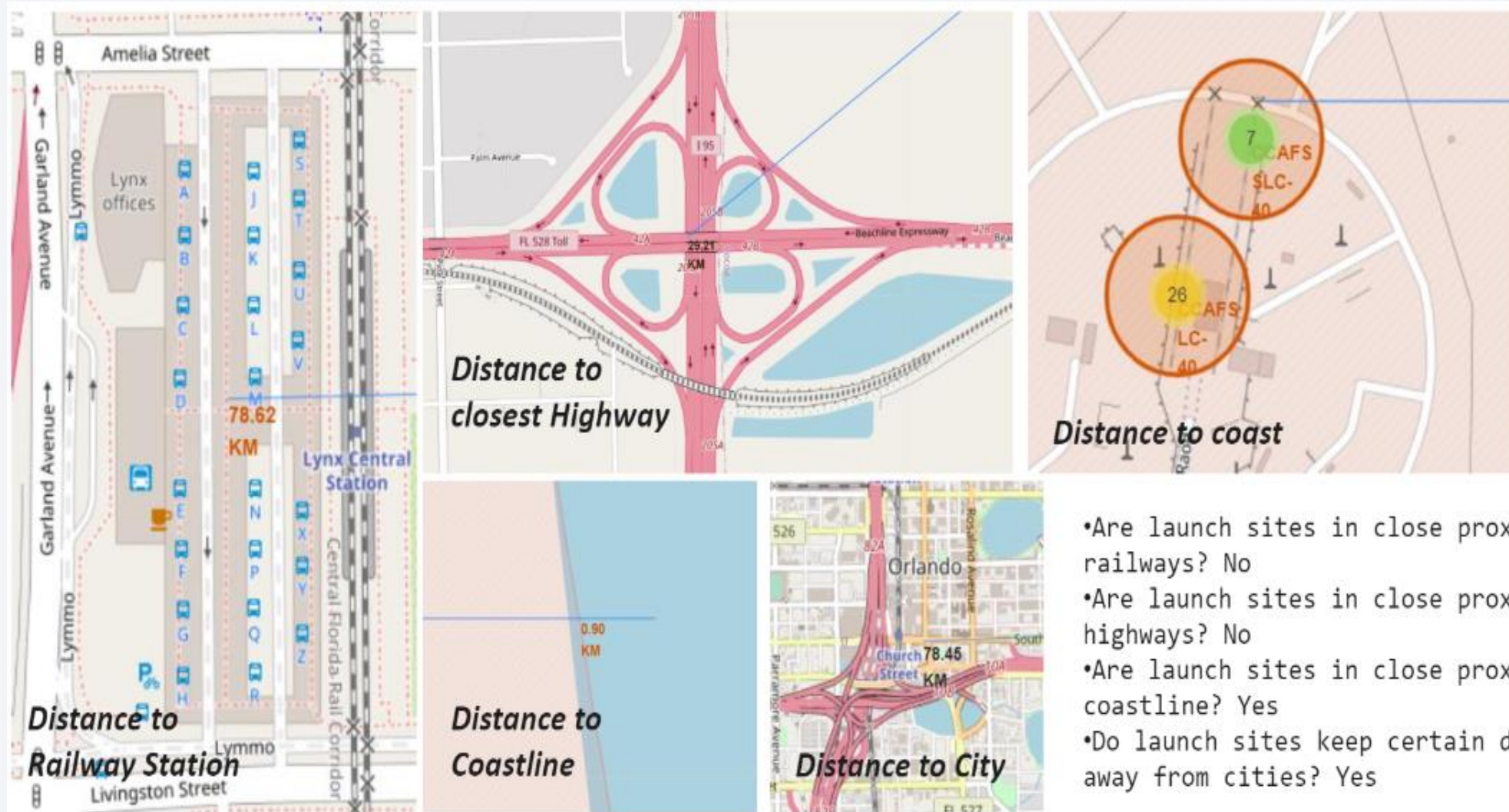


## <Folium Map Screenshot 2>





# <Folium Map Screenshot 3>



- Are launch sites in close proximity to railways? No
- Are launch sites in close proximity to highways? No
- Are launch sites in close proximity to coastline? Yes
- Do launch sites keep certain distance away from cities? Yes



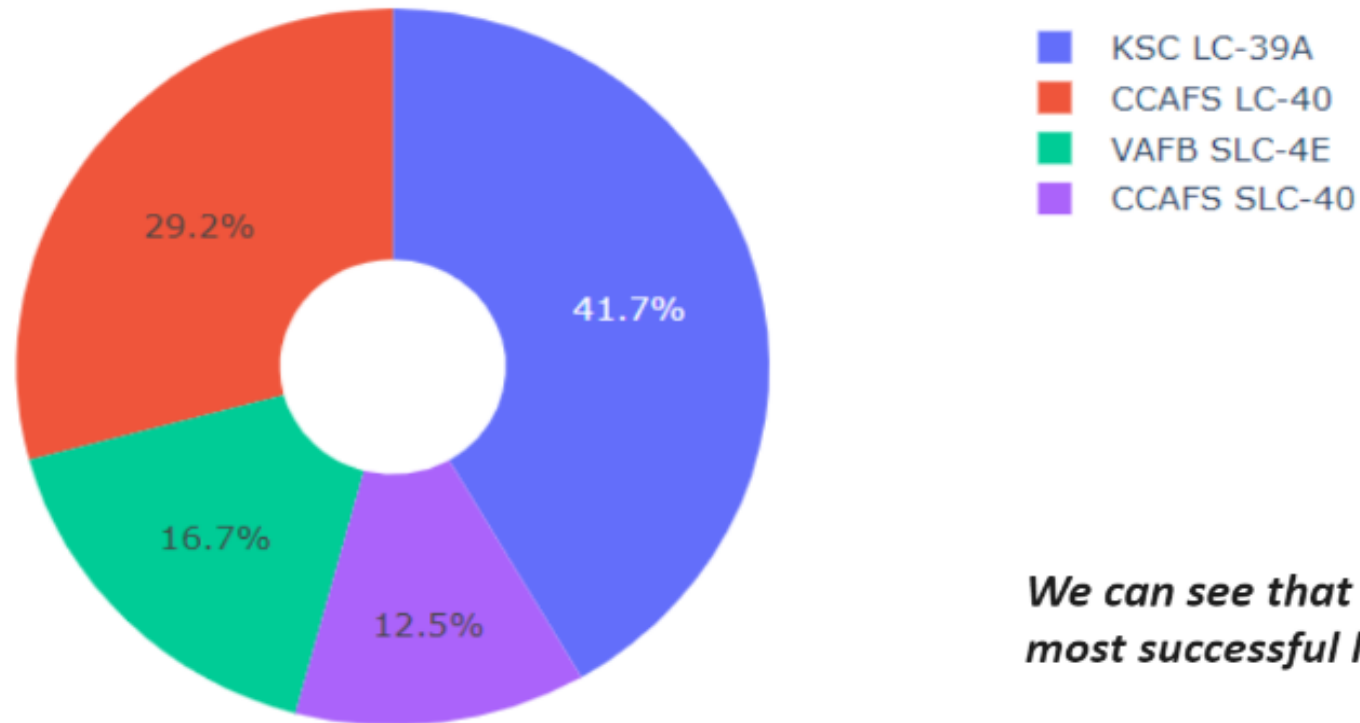


Section 4

# Build a Dashboard with Plotly Dash

# Pie Chart achievements by launch site

Total Success Launches By all sites

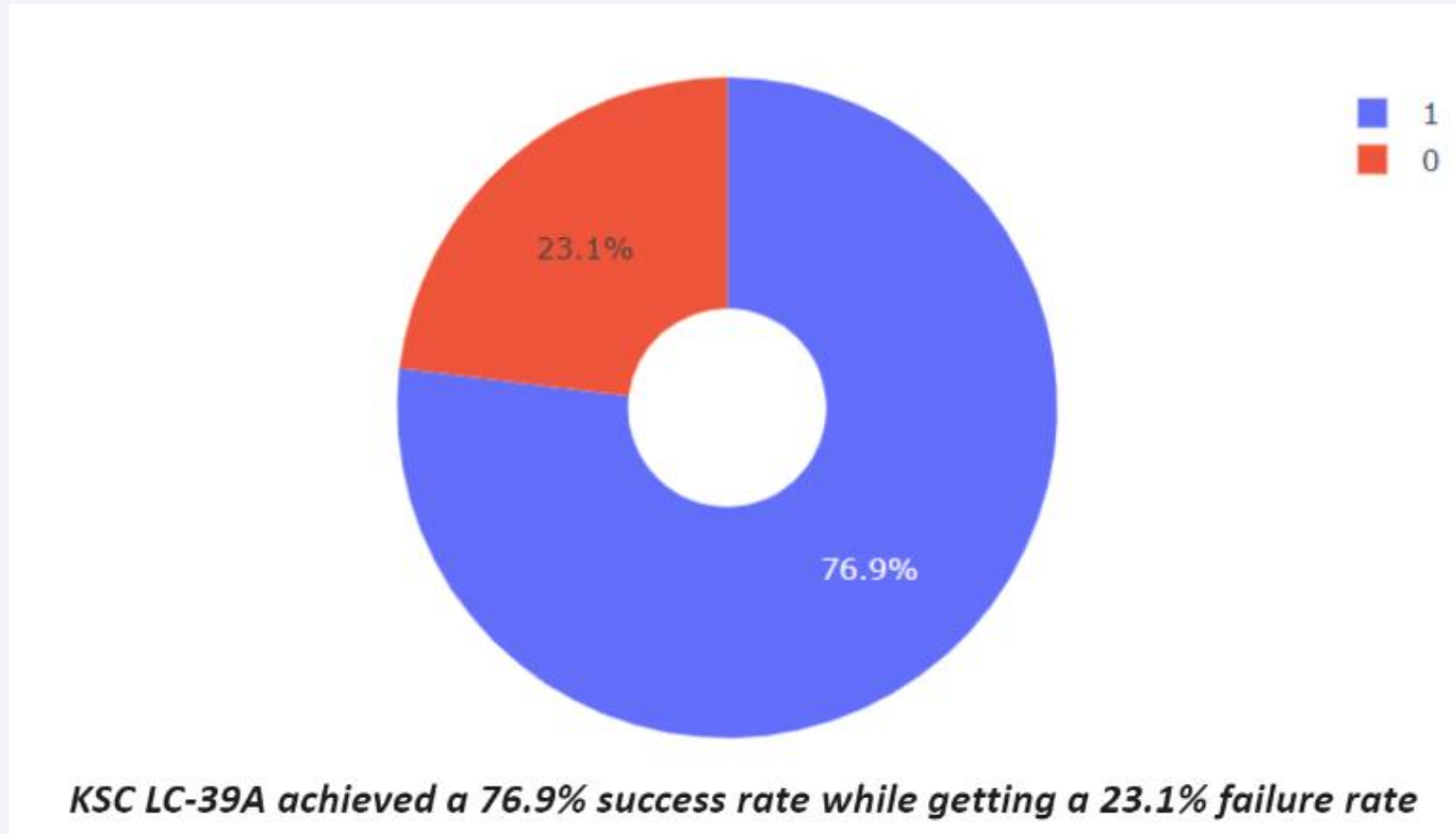


*We can see that KSC LC-39A had the most successful launches from all the sites*

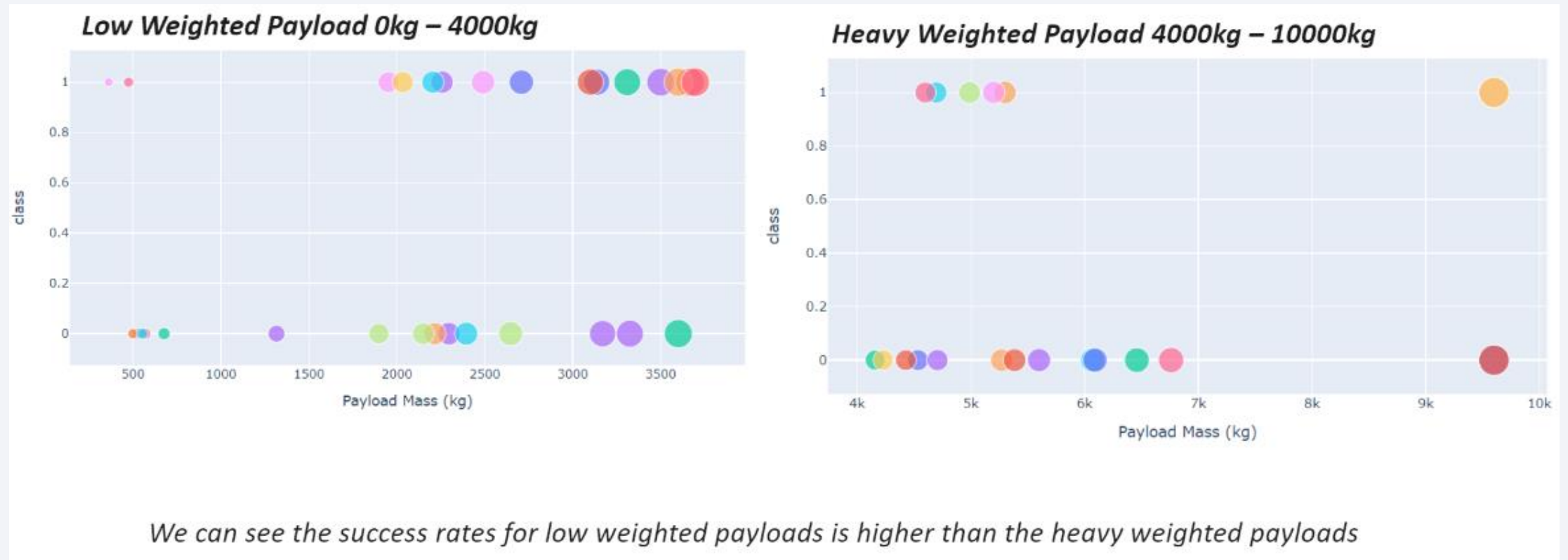


# Highest Launch Success Ratio Pie Chart

---



# Payload / Launch Outcome



Section 5

# Predictive Analysis (Classification)

# Classification Accuracy

---

- The decision tree classifier is the model with the highest classification accuracy

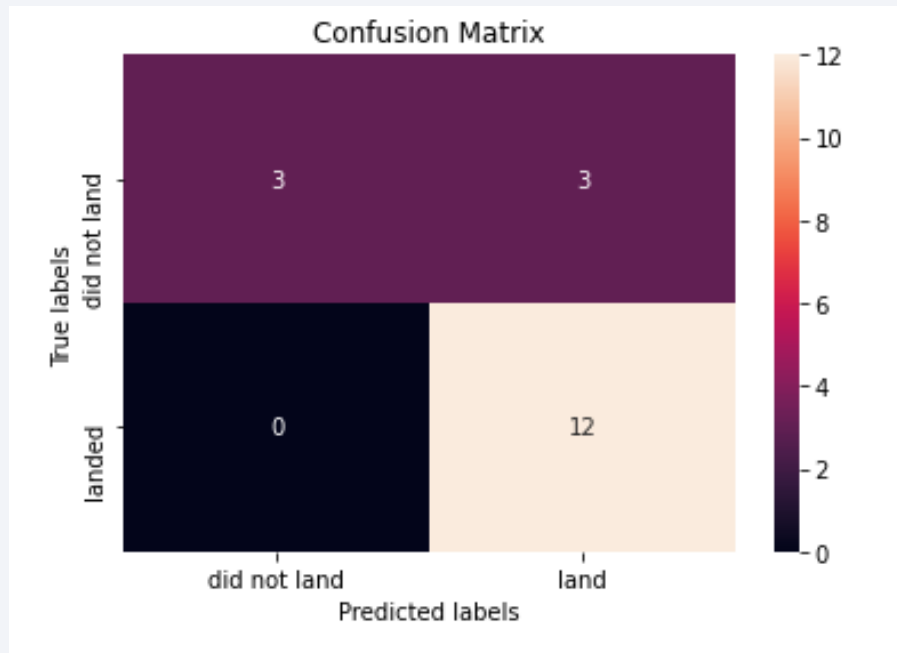
```
models = {'KNeighbors': knn_cv.best_score_,
          'DecisionTree': tree_cv.best_score_,
          'LogisticRegression': logreg_cv.best_score_,
          'SupportVector': svm_cv.best_score_}

bestalgorithm = max(models, key=models.get)
print('Best model is', bestalgorithm, 'with a score of', models[bestalgorithm])
if bestalgorithm == 'DecisionTree':
    print('Best params is :', tree_cv.best_params_)
if bestalgorithm == 'KNeighbors':
    print('Best params is :', knn_cv.best_params_)
if bestalgorithm == 'LogisticRegression':
    print('Best params is :', logreg_cv.best_params_)
if bestalgorithm == 'SupportVector':
    print('Best params is :', svm_cv.best_params_)
```

Best model is DecisionTree with a score of 0.8732142857142856

Best params is : {'criterion': 'gini', 'max\_depth': 6, 'max\_features': 'auto', 'min\_samples\_leaf': 2, 'min\_samples\_split': 5, 'splitter': 'random'}

# Confusion Matrix



- The confusion matrix for the decision tree classifier shows that the classifier can distinguish between the different classes. The major problem is the false positives .i.e., unsuccessful landing marked as successful landing by the classifier.

# Conclusions

---

- The larger quantity of flights 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.

# Appendix

---

- Include any relevant assets like Python code snippets, SQL queries, charts, Notebook outputs, or data sets that you may have created during this project



Thank you!

