

Applied Data Science Capstone Project

Alif Putra Dewantara

31 December 2022



Outline



EXECUTIVE
SUMMARY



INTRODUCTION



METHODOLOGY



RESULTS



CONCLUSION



APPENDIX



Executive Summary

- Summary of Methodologies
 - Data collection through API
 - Data collection with Web Scrapping
 - Exploratory Data Analysis (EDA) with SQL
 - Exploratory Data Analysis (EDA) with Data Visualization
 - Building a map visualization with Folium
 - Building a dashboard with Plotly Dash
 - Predictive Analysis (Classification)
- Summary of All Results
 - EDA results
 - Interactive analytics in screenshots
 - Predictive Analytics results from Machine Learning Lab

Introduction

- SpaceX is the most successful company in terms of commercial space travel. The company advertises its rocket launches, especially Falcon 9 as low as 62 million dollars when the other providers' cost is up to 165 million dollars. This cost is possible to make because SpaceX has revolutionary technologies in terms of its reusable rockets. As a data scientist at SpaceY, a startup company rivaling SpaceX, we need to create a machine learning pipeline to predict the landing outcome in the future. By that, we can also make more information in terms of bids against Space X.
- Question to be answered
 - How do the variables (payload mass, launch site, number of flight, orbits) affect the success of the first stage landing?
 - What is the rate of successful landings over the years?
 - What is the best algorithm that can be use for classification?



The background of the slide is a dense, abstract composition of three-dimensional numbers. The numbers, ranging from 0 to 9, are rendered in a light blue-grey color and are oriented in various directions, creating a sense of depth and movement. They appear to be floating or stacked, with some numbers being more prominent than others. The overall effect is a complex, data-driven visual texture.

Methodology

Section 1

Data Collection

- We obtain SpaceX data from two sources:
 - Open Source SpaceX REST API
 - (FlightNumber, Date, BoosterVersion, PayloadMass, Orbit, LaunchSite, Outcome, Flights, GridFins, Reused, Legs, LandingPad, Block, ReusedCount, Serial, Longitude, Latitude)
 - List of Falcon 9 and Falcon Heavy launches from Wikipedia (through web scrapping)
 - (Flight No., Launch site, Payload, PayloadMass, Orbit, Customer, Launch outcome, Version Booster, Booster landing, Date, Time)

[illegible]

Data Collection through SpaceX API



Get rocket launch data from SpaceX API



Use .json() to decode response and convert it to dataframe using .json_normalize()



Making data we already have into a dictionary



Filtering the dataframe to only include the data we needed



Performed data cleaning and filling the missing value by applying custom function



Exporting the data to .csv

Source:

[Applied-Data-Science-Capstone-IBM/1. jupyter-labs-spacex-data-collection-api.ipynb](https://github.com/alifdewantaraa/Applied-Data-Science-Capstone-IBM/blob/main/1.%20jupyter-labs-spacex-data-collection-api.ipynb) at main · alifdewantaraa/Applied-Data-Science-Capstone-IBM (github.com)



```
spacex_url="https://api.spacexdata.com/v4/launches/past"
```

```
response = requests.get(spacex_url)
```

```
# Use json_normalize meethod to convert the json result  
data = pd.json_normalize(response.json())
```

```
# Create a data from launch_dict  
df = pd.DataFrame.from_dict(launch_dict)
```

```
data_falcon9 = df[df['BoosterVersion']!='Falcon 1']
```

Now that we have removed some values we should reset the FlightNumber column

```
data_falcon9.loc[:, 'FlightNumber'] = list(range(1, data_falcon9.shape[0]+1))  
data_falcon9
```

```
# Calculate the mean value of PayloadMass column  
payloadmassavg = data_falcon9['PayloadMass'].mean()  
# Replace the np.nan values with its mean value  
data_falcon9['PayloadMass'].replace(np.nan, payloadmassavg, inplace=True)
```

```
data_falcon9.to_csv('dataset_part_1.csv', index=False)
```

```
# use requests.get() method with the provided static_url
# assign the response to a object
response = requests.get(static_url)
```

Create a `BeautifulSoup` object from the HTML `response`

```
# Use BeautifulSoup() to create a BeautifulSoup object from a response text cont
soup = BeautifulSoup(response.text, 'html')
```

Print the page title to verify if the `BeautifulSoup` object was created properly

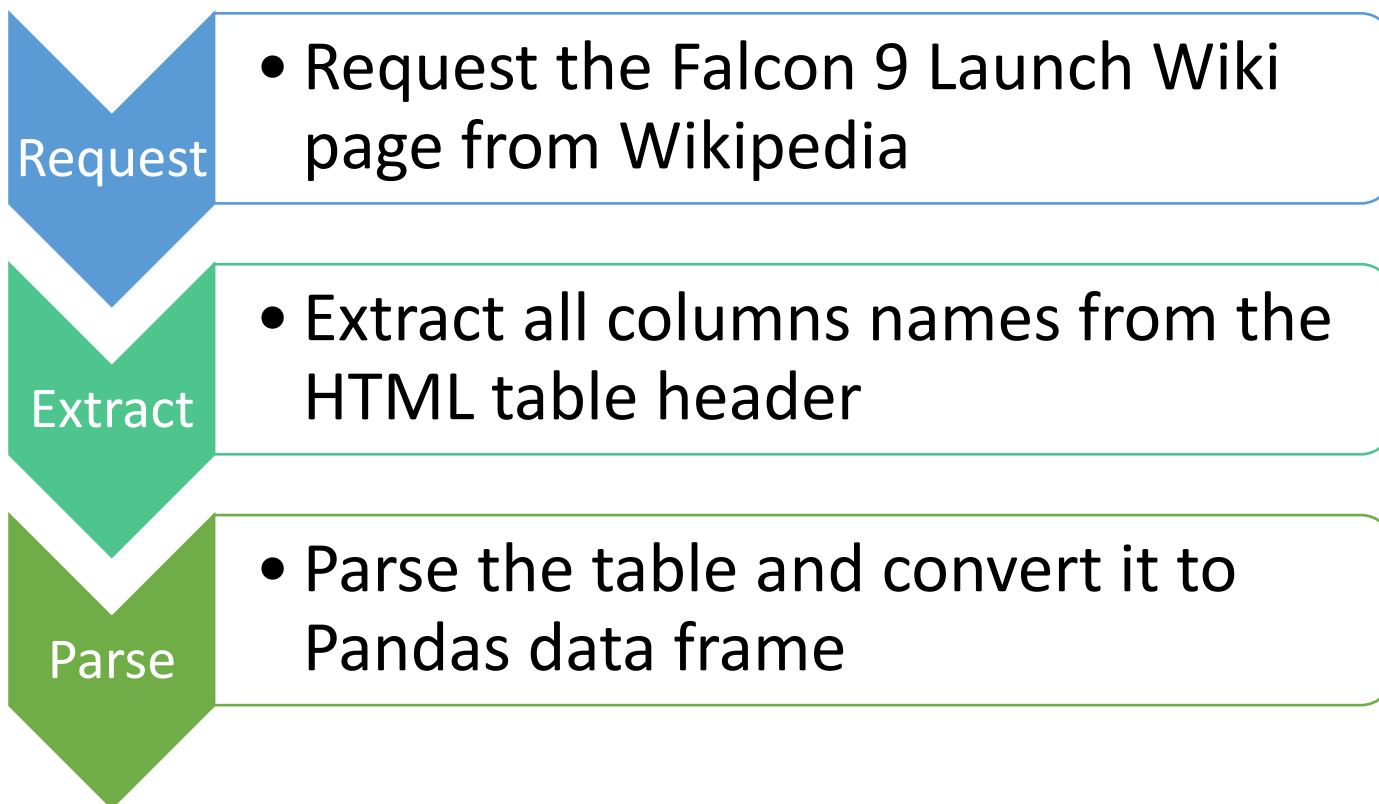
```
# Use soup.title attribute
soup.title
```

```
<title>List of Falcon 9 and Falcon Heavy launches - Wikipedia</title>
```

```
extracted_row = 0
#Extract each table
for table_number,table in enumerate(soup.find_all('table',"wikitable plainrowhea
# get table row
for rows in table.find_all("tr"):
    #check to see if first table heading is as number corresponding to launch
    if rows.th:
        if rows.th.string:
            flight_number=rows.th.string.strip()
            flag=flight_number.isdigit()
    else:
        flag=False
    #get table element
    row=rows.find_all('td')
    #if it is number save cells in a dictionary
    if flag:
```

```
df=pd.DataFrame(launch_dict)
```

Data Collection with Web Scrapping



Data Wrangling

Explore data

- Calculate the number of launches on each site
- Calculate the number and occurrence of each orbit and mission outcome per orbit

Create a landing outcome training label from 'Outcome' column

- True ASDS = drone ship landed succeeded
- None None = not attempted
- True RTLS = ground pad landed succeeded
- False ASDS = drone ship landing failed
- True Ocean = ocean landing succeeded
- None ASDS = unable to be attempted due to launch failure
- False Ocean = ocean landing failed
- False RTLS = ground pad landing failed

```
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: Outcome, dtype: int64
```

Exploratory Data Analysis (EDA) with SQL

The following SQL queries were performed:

- Names of the unique launch sites in the space mission;
- Top 5 launch sites whose name begin 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 drone ship, their booster versions, and launch site names for in year 2015; and
- Rank of the count of landing outcomes (such as Failure (drone ship) or Success (ground pad)) between the date 2010-06-04 and 2017-03-20.

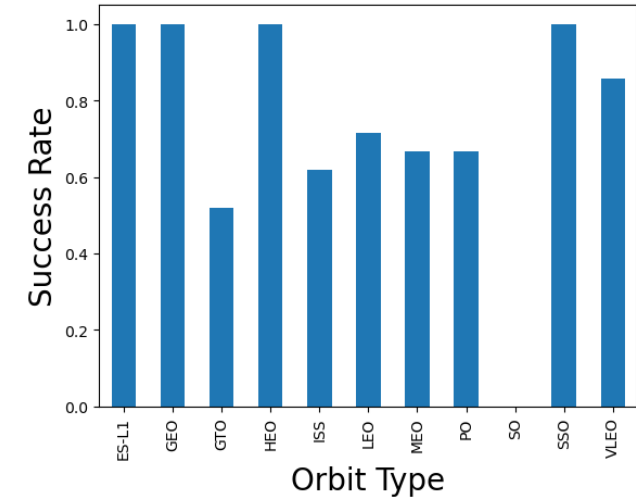


Exploratory Data Analysis (EDA) with Data Visualization

To explore data, scatterplots and barplots 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

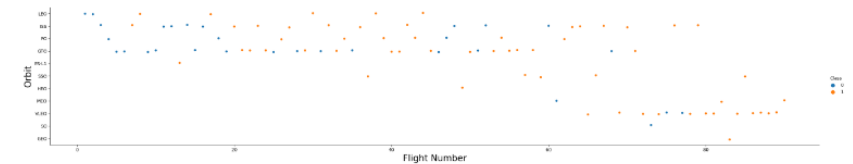
```
# HINT use groupby method on Orbit column and get the mean of Class column
df_orbit = df.groupby('Orbit')['Class'].mean()
df_orbit.plot(kind='bar')
plt.xlabel('Orbit Type', fontsize=20)
plt.ylabel('Success Rate', fontsize=20)
plt.show()
```



TASK 4: Visualize the relationship between FlightNumber and Orbit type

For each orbit, we want to see if there is any relationship between FlightNumber and Orbit type.

```
# Plot a scatter point chart with x axis to be FlightNumber and y axis to be the Orbit, and hue to be the
sns.catplot(x='FlightNumber', y='Orbit', hue='Class', data=df, aspect = 5)
plt.xlabel('Flight Number', fontsize = 20)
plt.ylabel('Orbit', fontsize = 20)
plt.show()
```



Building a map visualization 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.

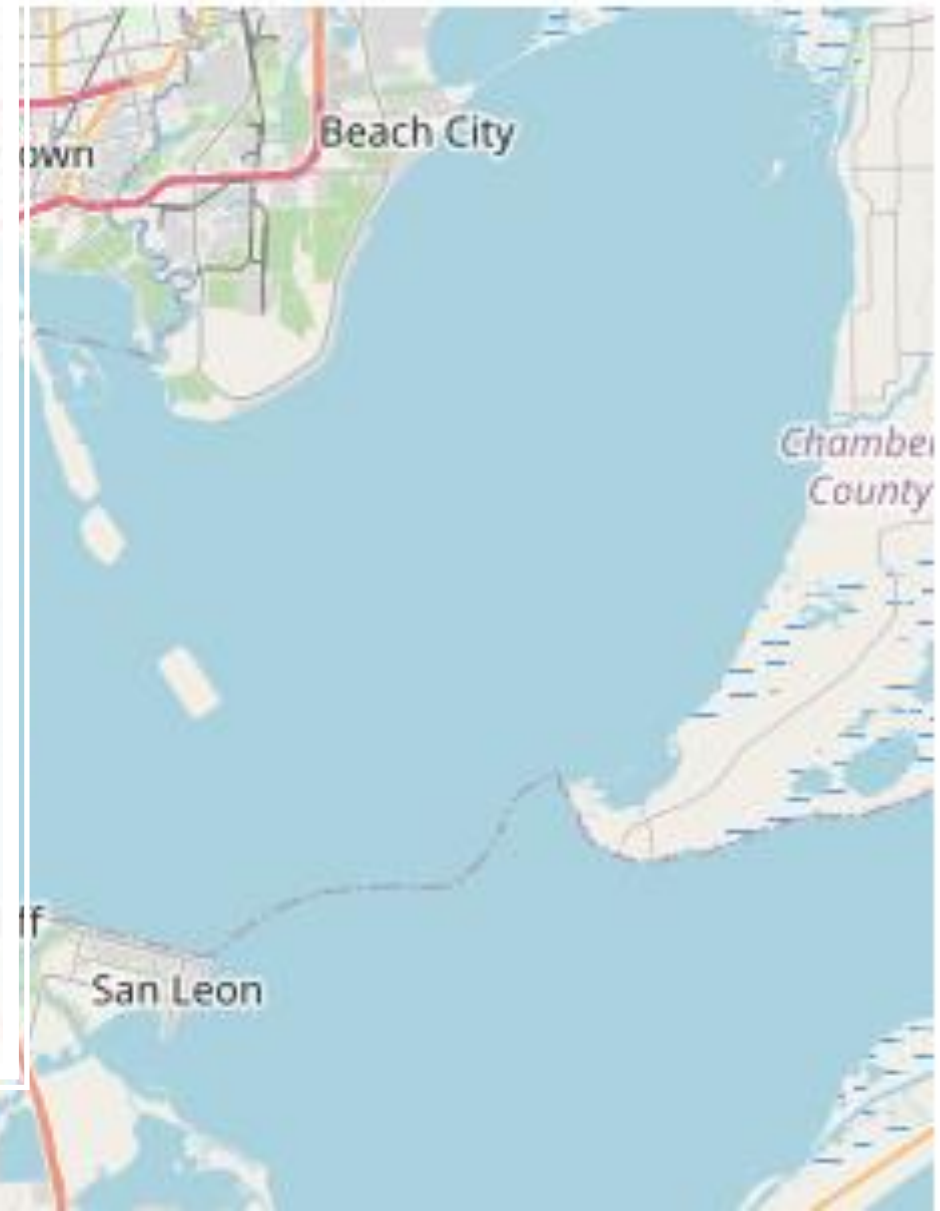
Colored Markers of the launch outcomes for each Launch Site:

- Added colored 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:

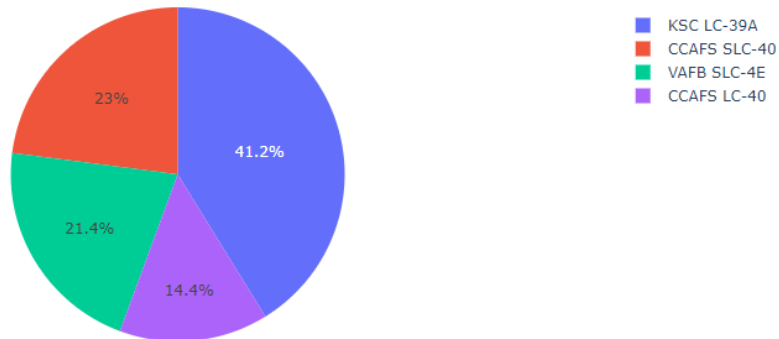
- Added colored Lines to show distances between the Launch Site KSC LC-39A (as an example) and its proximities like Railway, Highway, Coastline and Closest City

Source: [Applied-Data-Science-Capstone-IBM/6. lab jupyter launch site location.ipynb at main · alifdewantaraa/Applied-Data-Science-Capstone-IBM \(github.com\)](#)



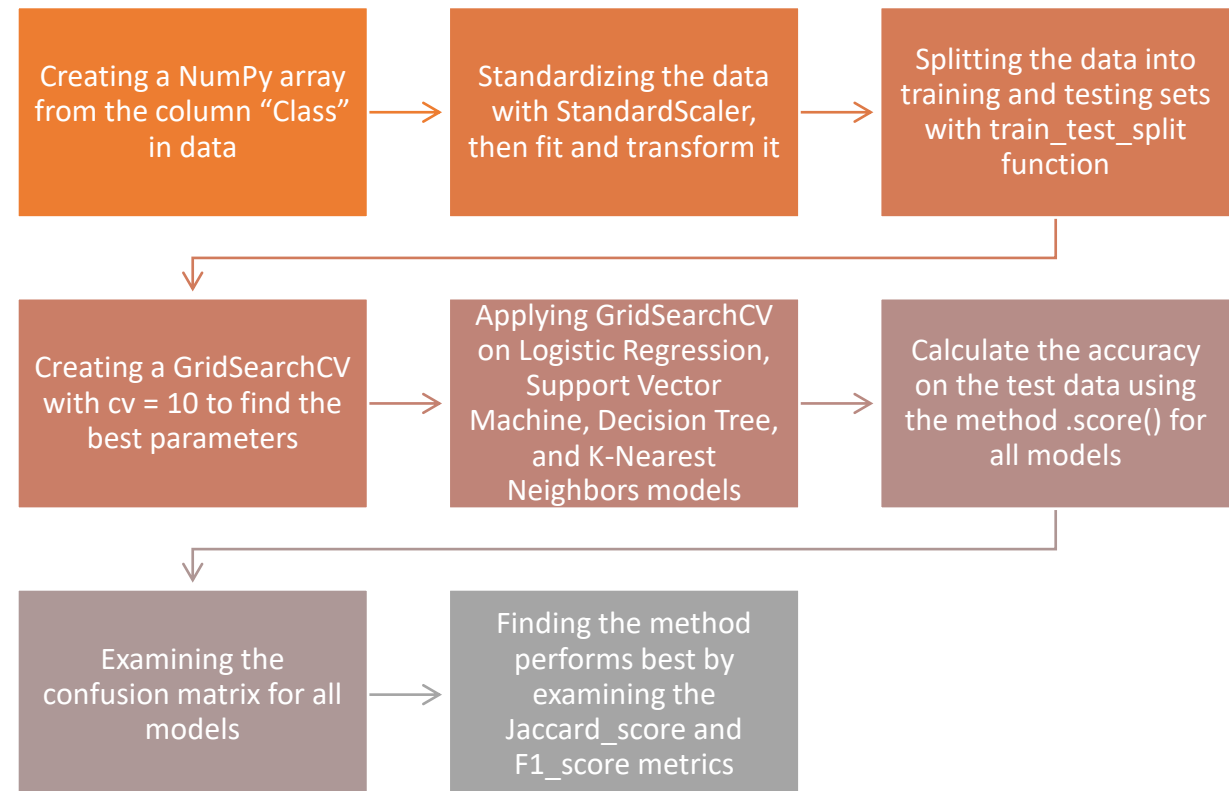
Building a Dashboard with Plotly Dash

Total Success Launches by Site



- 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.

Predictive Analysis (Classification)



Source: [Applied-Data-Science-Capstone-IBM/8. SpaceX Machine Learning Prediction Part 5.ipynb at main · alifdewantaraa/Applied-Data-Science-Capstone-IBM \(github.com\)](#)

Results

The results will be categorized to main results:

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

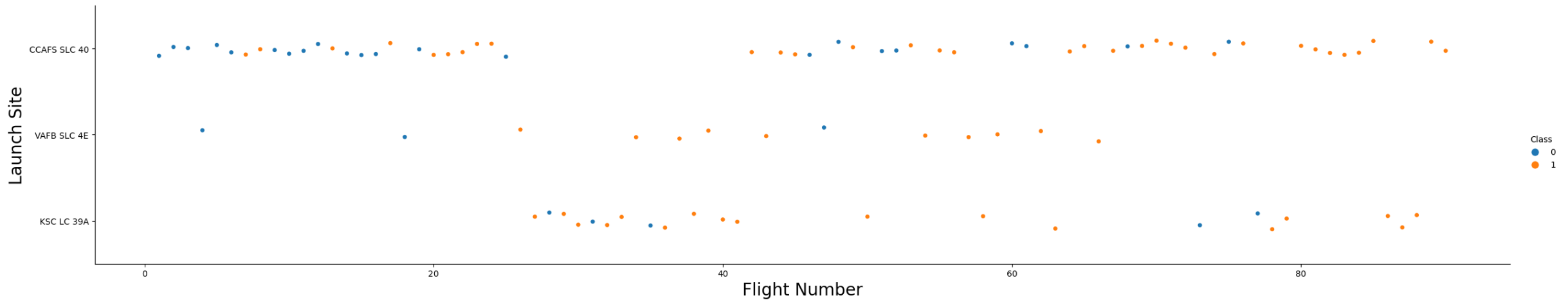


Insights drawn from Exploratory Data Analysis

Section 2

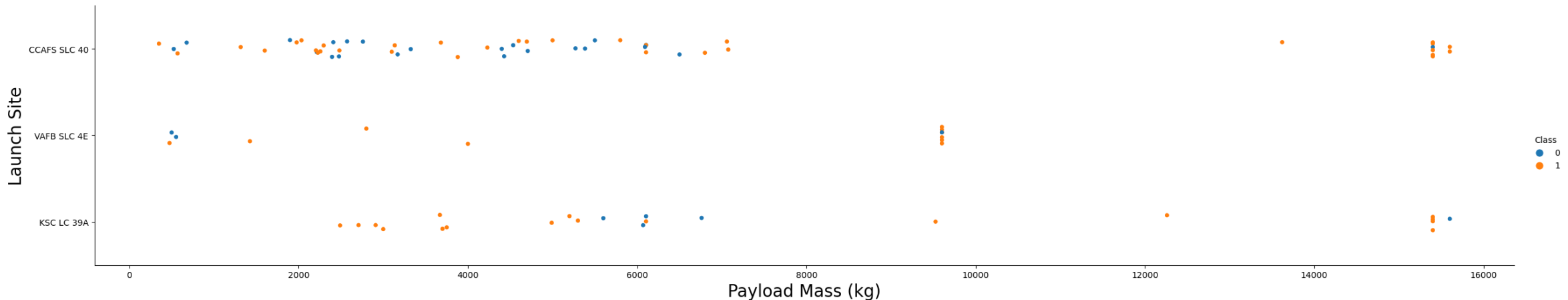
Flight Number vs. Launch Site

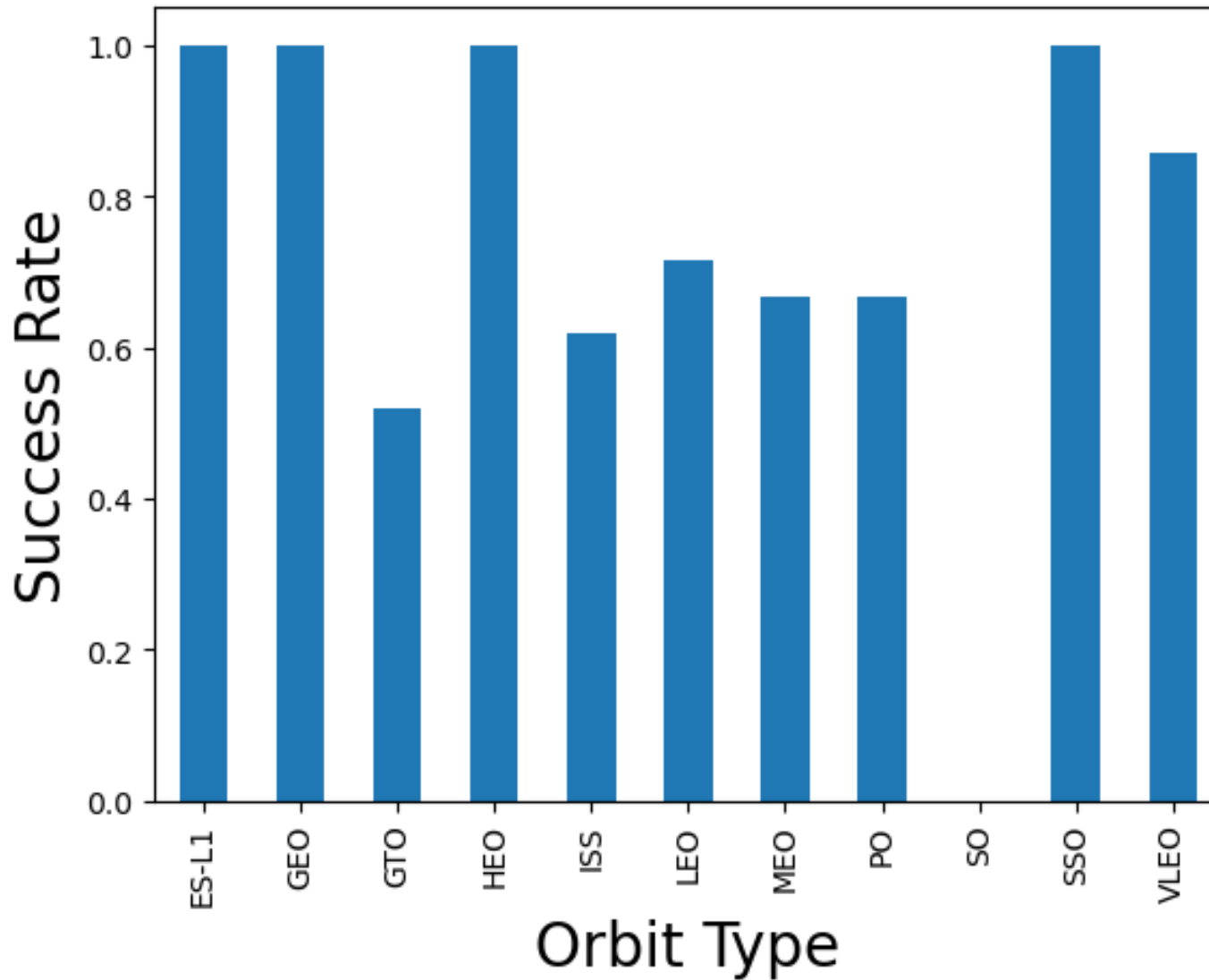
- This plot above explains that it is possible to verify that the best launch site nowadays is CCAF5 SLC 40
- The earliest flights all failed while the latest flight all succeeded
- VAFB SLC 4E and KSC LC 39A have higher success rates
- It also explains that each new launch has higher rate of success



Payload 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.





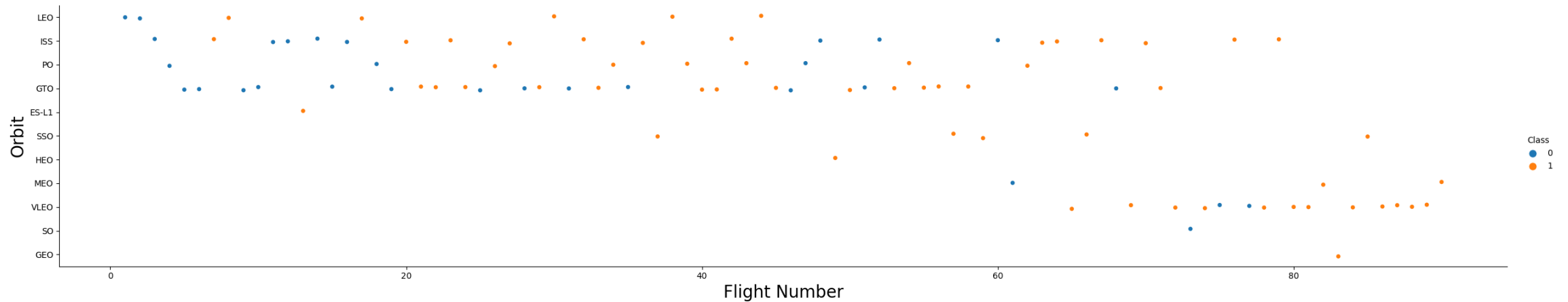
Success Rate vs. Orbit Type

- 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, P

However, deeper analysis show that some of this orbits has only 1 occurrence such as GEO, SO, HEO and ES-L1 which mean this data need more dataset to see pattern or trend before we draw any conclusion.

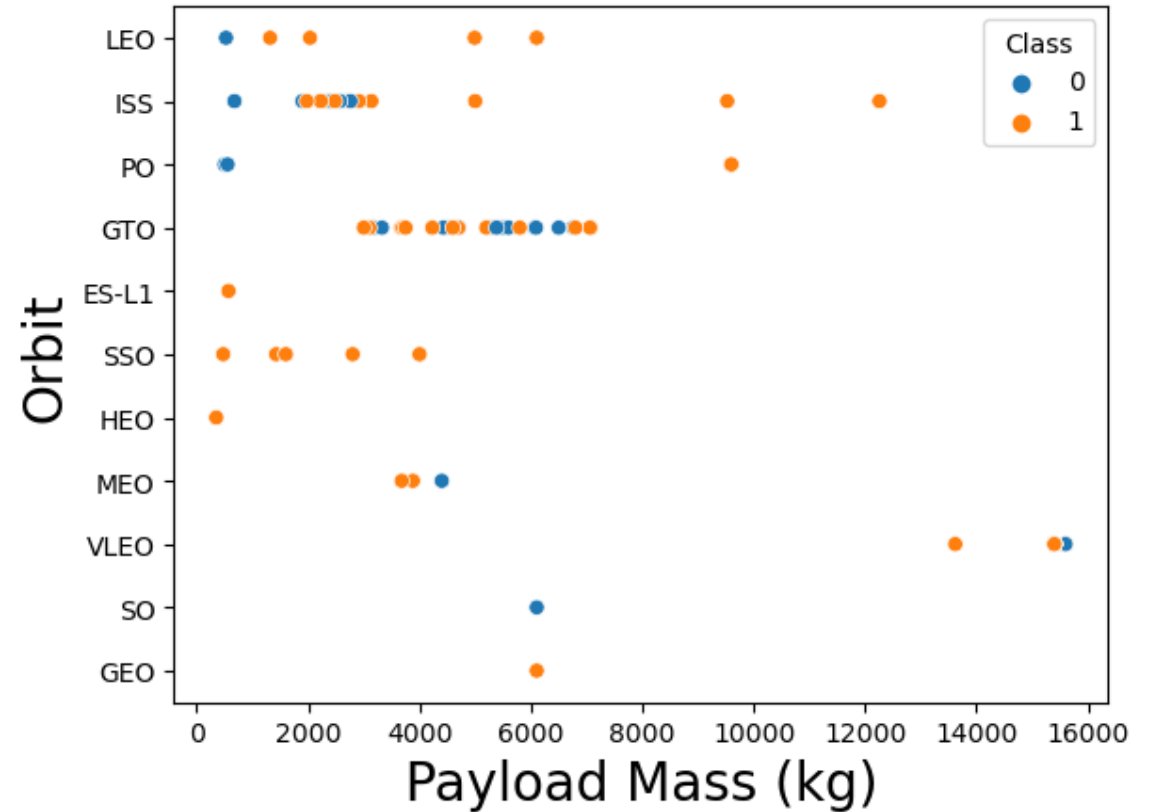
Flight Number vs. Orbit type

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.



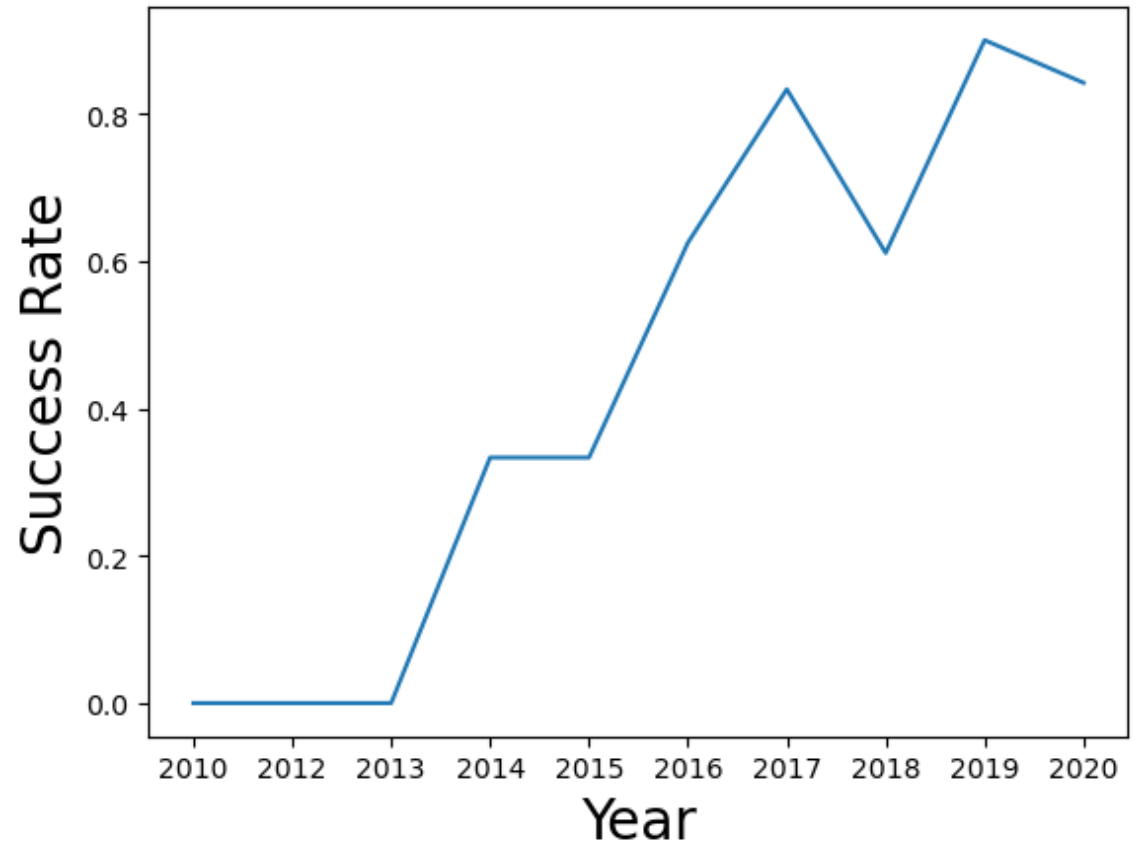
Payload Mass vs. Orbit Type

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



Launch Success Yearly Trend

This figure clearly depicts an increasing trend from the year 2013 until 2020. If this trend continues for the next year onward, the success rate will steadily increase until reaching 100% success rate.



All Launch Site Names

To show only unique launch sites from SpaceX data, we used keyword **DISTINCT**

```
%sql SELECT DISTINCT launch_site FROM SPACEX
```

```
* ibm_db_sa://hwz01292:***@21fecfd8-47b7-4937-840d-d791d021866  
0.bs2io90108kqb1od81cg.databases.appdomain.cloud:31864/blddb  
Done.
```

launch_site

CCAFS LC-40

CCAFS SLC-40

KSC LC-39A

VAFB SLC-4E

Launch Site Names Begin with 'CCA'

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

Display 5 records where launch sites begin with the string 'CCA'

```
%sql SELECT * FROM SpaceX WHERE launch_site LIKE 'CCA%' LIMIT 5
```

```
* ibm_db_sa://hwz01292:***@21fecfd8-47b7-4937-840d-d791d0218660.bs2io90108kqb1od8lcg.databases.appdomain.cloud:31864/bludb  
Done.
```

DATE	time_utc_	booster_version	launch_site	payload	payload_mass_kg_	orbit	customer	mission_outcome	landing_outcome
2010-04-06	18:45:00	F9 v1.0 B0003	CCAFS LC-40	Dragon Spacecraft Qualification Unit	0	LEO	SpaceX	Success	Failure (parachute)
2010-08-12	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)
2012-08-10	00:35:00	F9 v1.0 B0006	CCAFS LC-40	SpaceX CRS-1	500	LEO (ISS)	NASA (CRS)	Success	No attempt
2013-01-03	15:10:00	F9 v1.0 B0007	CCAFS LC-40	SpaceX CRS-2	677	LEO (ISS)	NASA (CRS)	Success	No attempt
2013-03-12	22:41:00	F9 v1.1	CCAFS LC-40	SES-8	3170	GTO	SES	Success	No attempt

Total payload mass

We calculated the total payload carried by boosters from NASA as **22007** using the query below

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

```
%sql SELECT SUM(payload_mass__kg_) AS total FROM SpaceX WHERE customer = 'NASA (CRS)'
```

```
* ibm_db_sa://hwz01292:***@21fecfd8-47b7-4937-840d-d791d0218660.bs2io90108kqb1od81cg.databases.appdomain.cloud:31864/bludb
```

Done.

total

22007

Average Payload Mass by F9 v1.1

We calculated the average payload mass carried by booster version F9 v1.1 as **3678**

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

```
%sql SELECT AVG(payload_mass__kg_) AS average FROM SpaceX WHERE booster_version = 'F9 v1.1'
```

```
* ibm_db_sa://hwz01292:***@21fecfd8-47b7-4937-840d-d791d0218660.bs2io90l08kqb1od8lcg.databases.appdomain.cloud:31864/bludb
```

```
Done.
```

```
average
```

```
3676
```

First Successful Ground Landing Date

- We use the *min()* function to find the result. We observed that the dates of the first successful landing outcome on drone ship pad was **05th July 2016**.

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

Hint: Use min function

```
%sql SELECT MIN(DATE) FROM SpaceX WHERE landing__outcome = 'Success (drone ship)'
```

```
* ibm_db_sa://hwz01292:***@21fecfd8-47b7-4937-840d-d791d0218660.bs2io90108kqb1od8lcg.databases.appdomain.cloud:31864/bludb  
Done.
```

1

2016-06-05

Successful drone ship
landing with payload
between 4000 and 6000

- We used the WHERE clause to filter for boosters which have successfully landed on drone ship and applied the AND condition to determine successful landing with payload mass between 4000 and 6000

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

```
%sql SELECT booster_version FROM SpaceX WHERE landing__outcome = 'Success (drone ship)' AND payload_mass__kg_ BETWEEN 4000 AND 6000
```

```
* ibm_db_sa://hwz01292:***@21fecfd8-47b7-4937-840d-d791d0218660.bs2io90l08kqb1od8l1cg.databases.appdomain.cloud:31864/bludb  
Done.
```

booster_version

F9 FT B1022

F9 FT B1031.2

List the total number of successful and failure mission outcomes

```
%sql SELECT mission_outcome, COUNT(mission_outcome) FROM SpaceX GROUP BY mission_outcome
```

```
* ibm_db_sa://hwz01292:***@21fecfd8-47b7-4937-840d-d791d0218660.bs2io90l08kqb1od8lcg.databases.appdomain.cloud:31864/bludb
```

Done.

mission_outcome	2
-----------------	---

Success	44
---------	----

Success (payload status unclear)	1
----------------------------------	---

Total number of
successful and
failure mission
outcomes



List the names of the booster_versions which have carried the maximum payload mass. Use a subquery

```
%sql SELECT booster_version FROM SpaceX WHERE payload_mass__kg_ = (SELECT MAX(payload_mass__kg_) FROM SpaceX)
```

```
* ibm_db_sa://hwz01292:***@21fecfd8-47b7-4937-840d-d791d0218660.bs2io90l08kqb1od8l1cg.databases.appdomain.cloud:31864/bludb  
Done.
```

booster_version

F9 B5 B1048.4

F9 B5 B1049.4

F9 B5 B1049.5

F9 B5 B1060.2

F9 B5 B1058.3

Boosters carried
maximum
payload



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

```
%sql SELECT landing__outcome, booster_version, launch_site FROM SpaceX WHERE landing__outcome = 'Failure (drone ship)' AND YEAR(DATE) = 2015
```

```
* ibm_db_sa://hwz01292:***@21fecfd8-47b7-4937-840d-d791d0218660.bs2io90l08kqb1od8lcg.databases.appdomain.cloud:31864/bludb  
Done.
```

landing__outcome	booster_version	launch_site
------------------	-----------------	-------------

Failure (drone ship)	F9 v1.1 B1012	CCAFS LC-40
----------------------	---------------	-------------

2015 launch records



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

```
%sql SELECT landing__outcome, COUNT(landing__outcome) AS count FROM SpaceX WHERE DATE BETWEEN '2010-06-04' AND '2017-03-20' GROUP BY landing__outcome
```

```
* ibm_db_sa://hwz01292:***@21fecfd8-47b7-4937-840d-d791d0218660.bs2io90l08kqb1od8lcg.databases.appdomain.cloud:31864/bludb  
Done.
```

landing__outcome	COUNT
No attempt	7
Failure (drone ship)	2
Success (drone ship)	2
Success (ground pad)	2
Controlled (ocean)	1
Failure (parachute)	1

Rank success count
between 2010-06-
04 and 2017-03-20

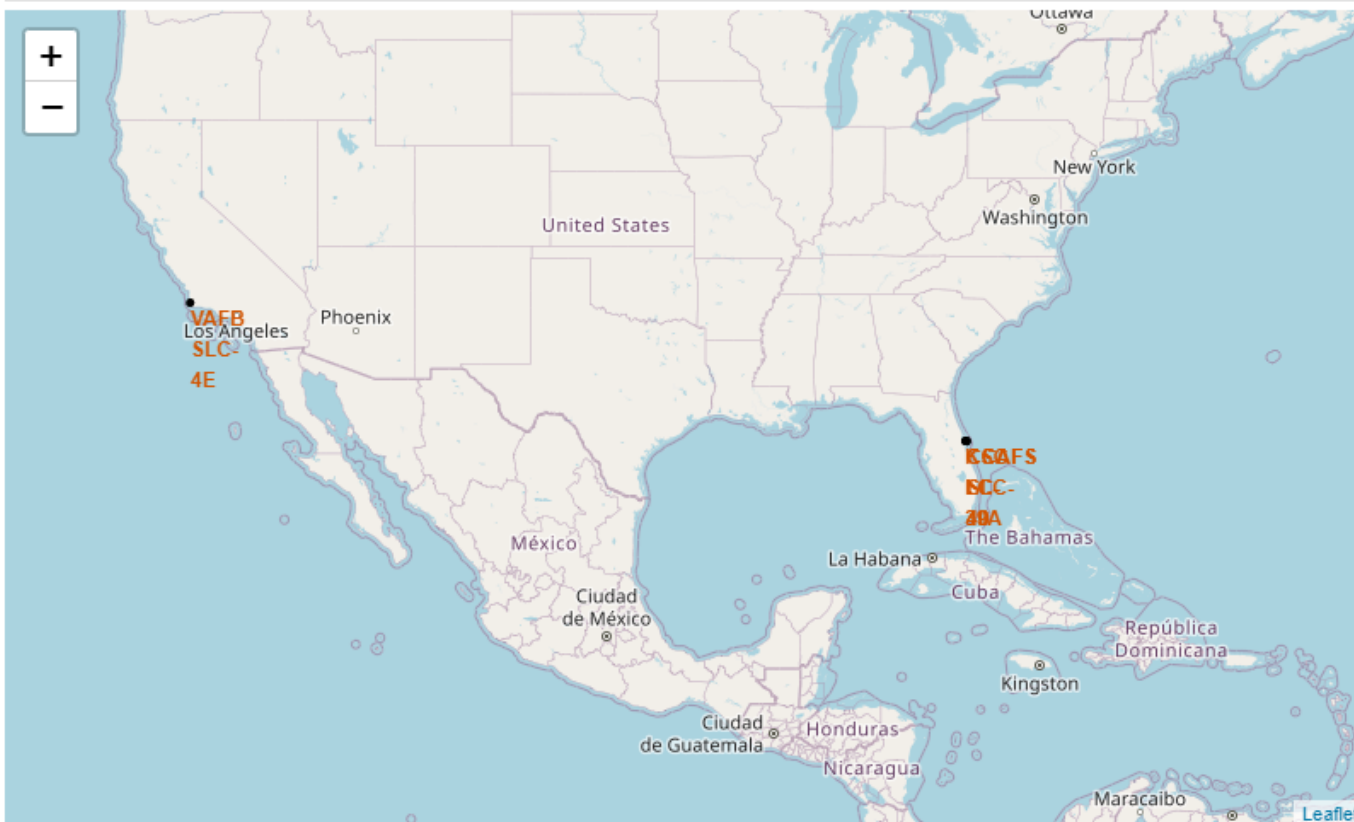




Launch Sites Proximities Analysis

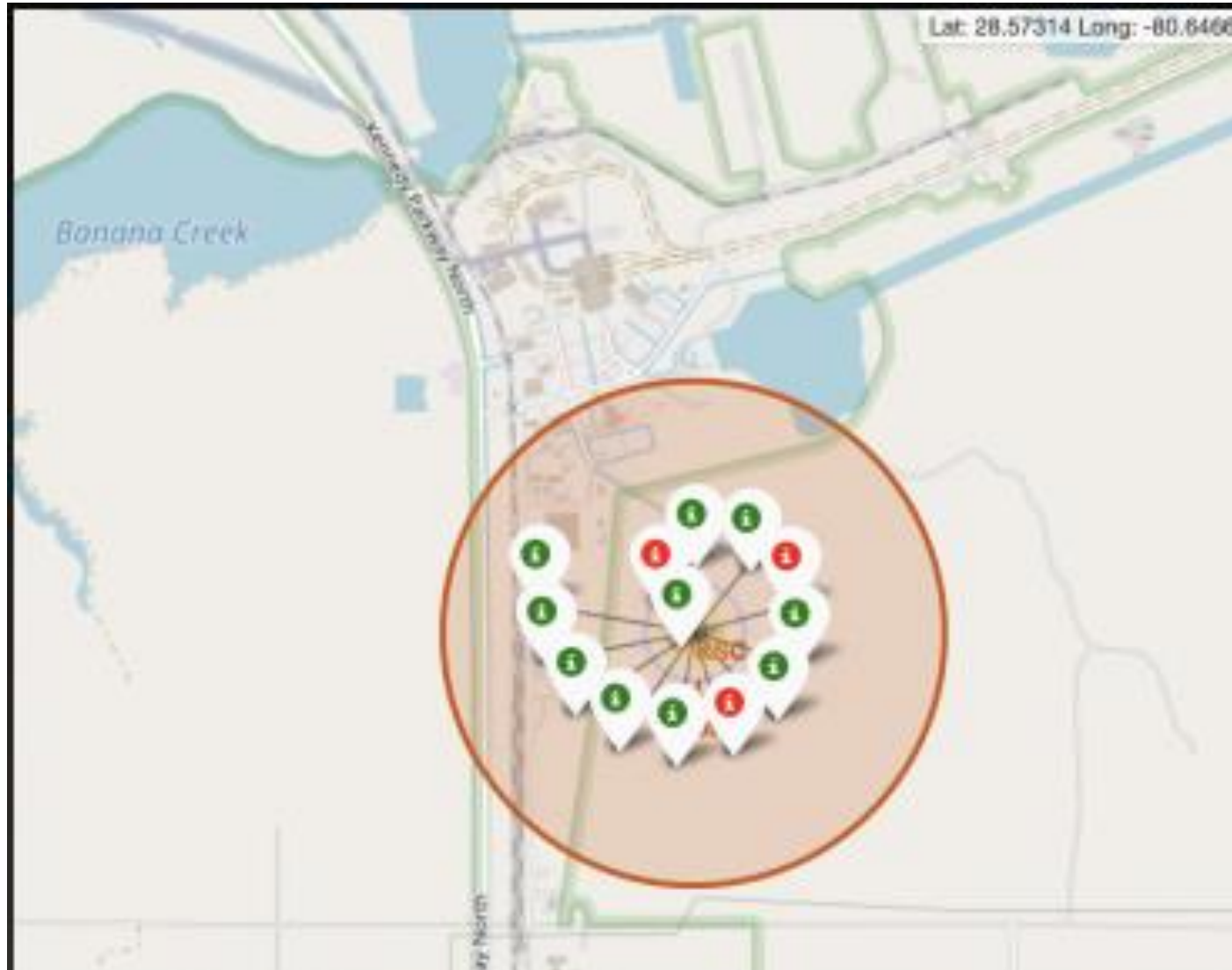
Section 3

```
# Initialize the map
site_map = folium.Map(location=nasa_coordinate, zoom_start=4)
# For each launch site, add a Circle object based on its coordinate (Lat, Long) values. In addition, add Launch site
for index, row in launch_sites_df.iterrows():
    coordinate = [row['Lat'], row['Long']]
    folium.Circle(coordinate, radius=1000, color='#000000', fill=True).add_child(folium.Popup(row['Launch Site']))
    folium.Marker(coordinate, icon=DivIcon(icon_size=(20,20),icon_anchor=(0,0), html='<div style="font-size: 1
site_map
```



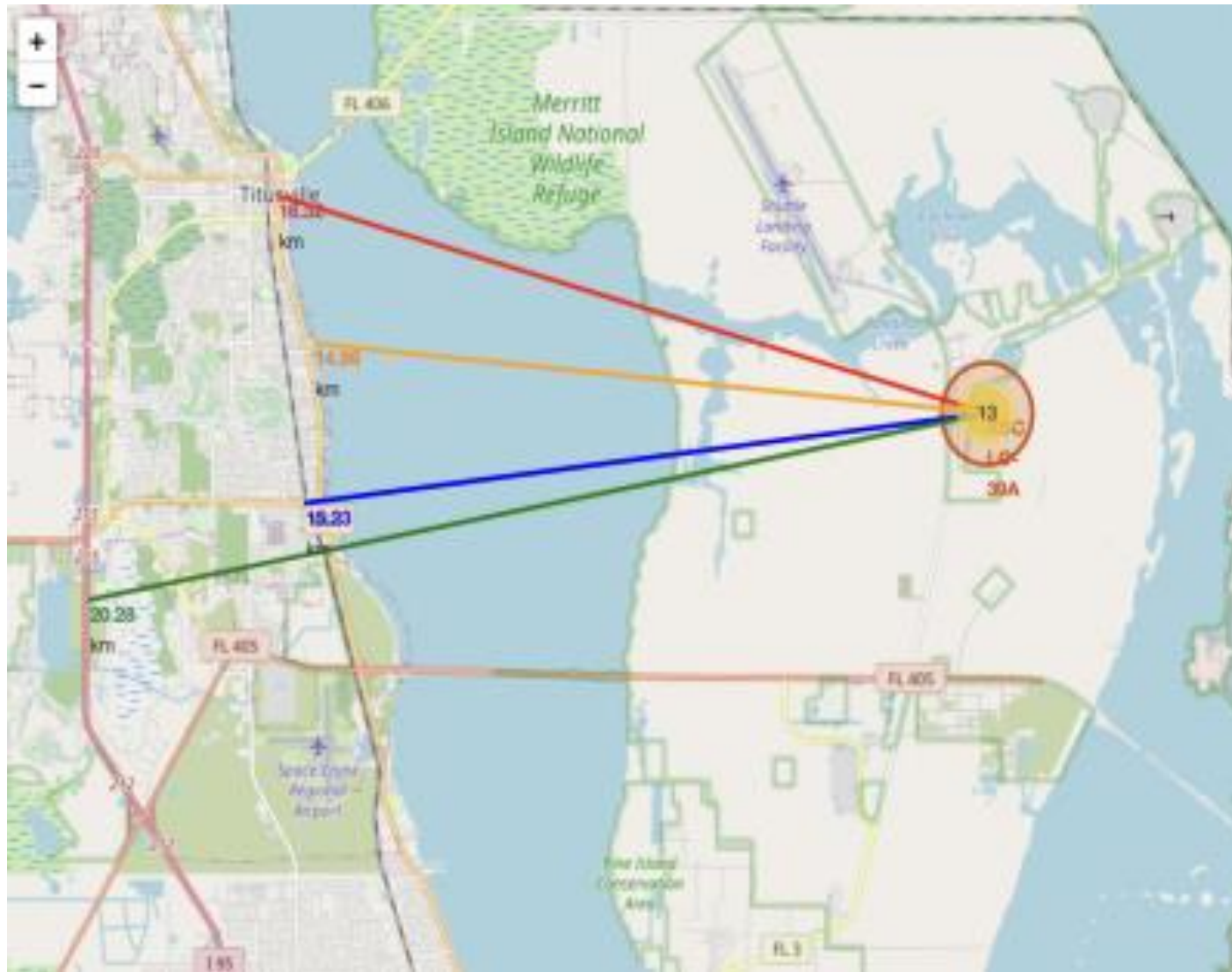
Location of all the Launch Sites

- Visualizing the launch sites on a map highlights the importance of launch site proximity to the coast and equator:



Colour-labeled launch records on the map

- 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

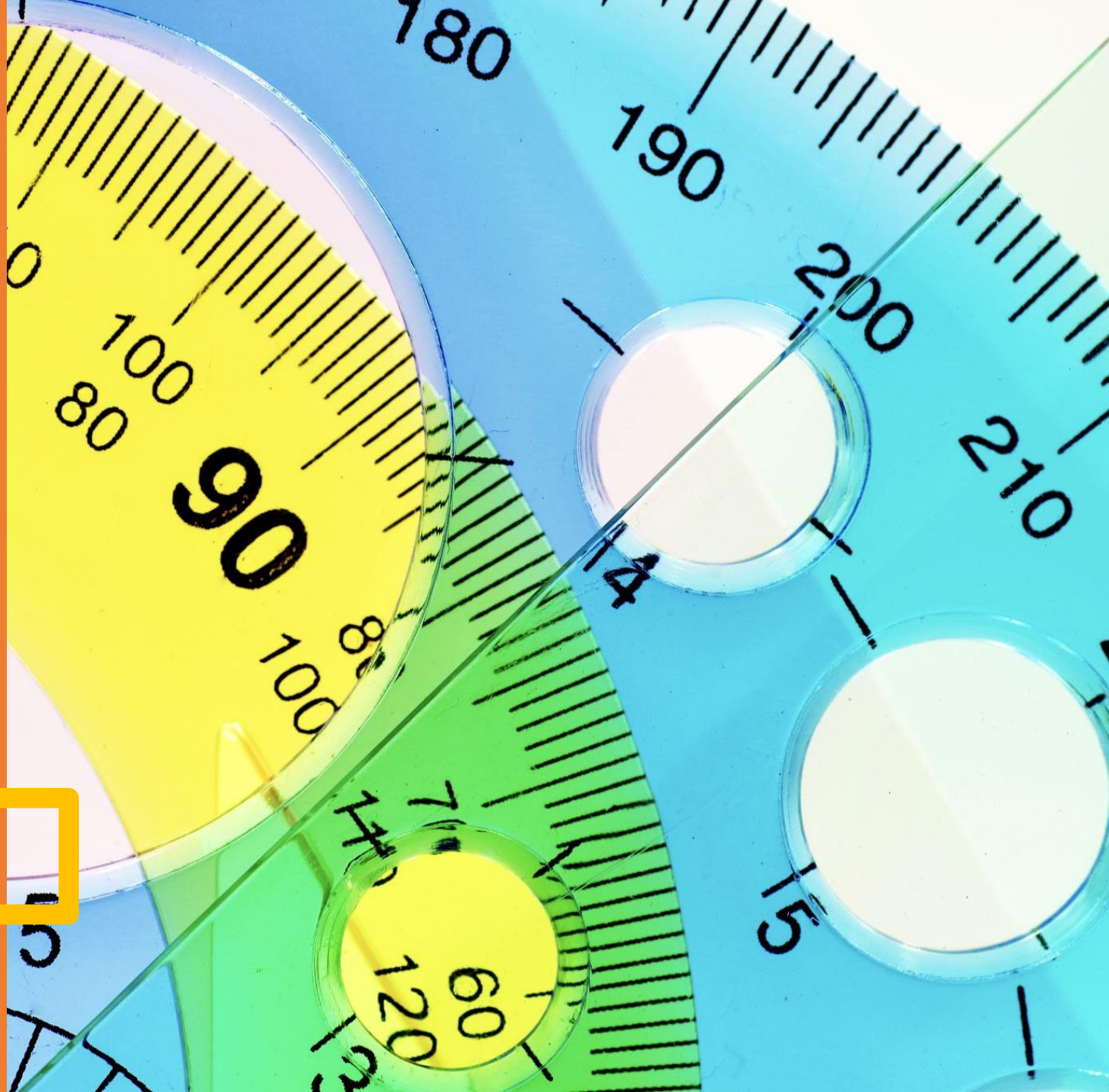


Distance from the launch site KSC LC-39A to its proximities

- From the visual analysis of the launch site KSC LC-39A we can clearly see that it is:
 - relatively close to railway (15.23 km)
 - relatively close to highway (20.28 km)
 - relatively close to coastline (14.99 km)
- Also, the launch site KSC LC-39A is relatively 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

Build a pie chart with Plotly Dash

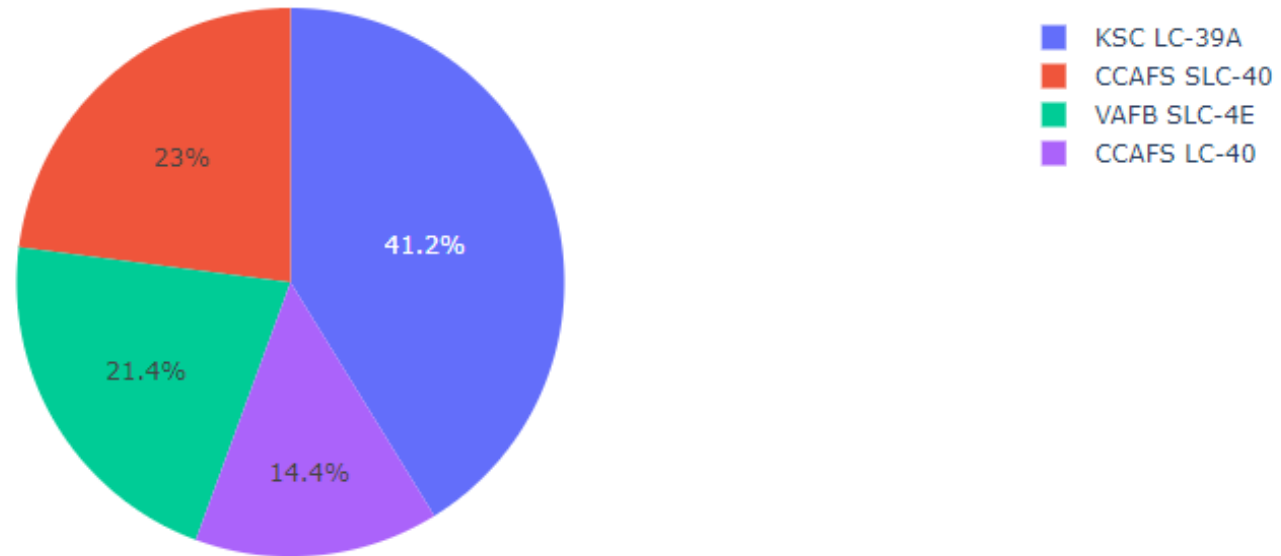
Section 4



The success percentage by each sites

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

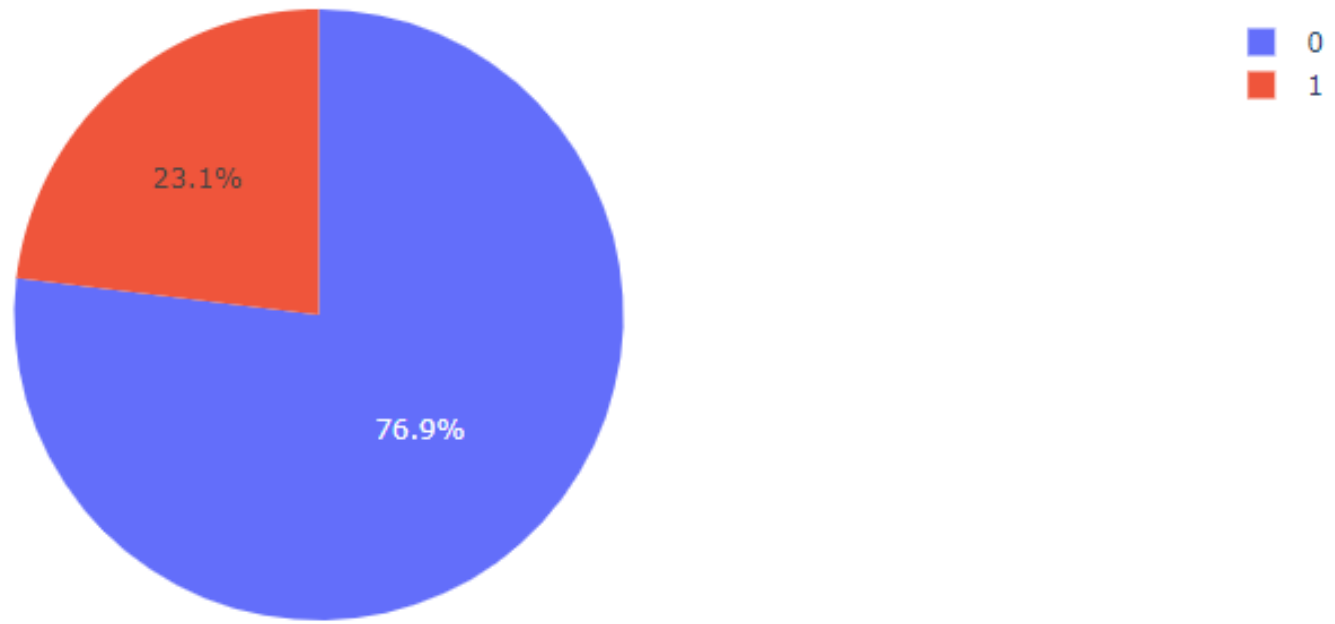
Total Success Launches by Site



Launch site with
highest launch
success ratio

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

Total Success Launches for Site KSC LC-39A



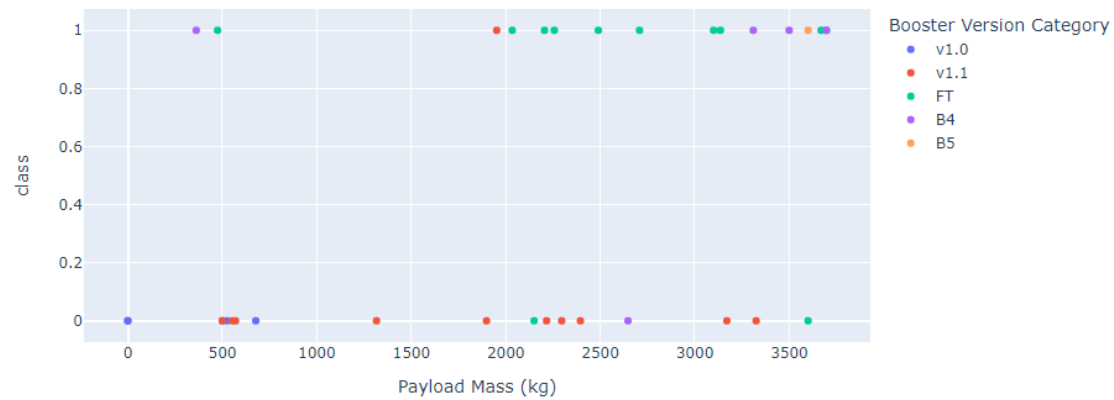
Payload vs Launch Outcome Scatter Plot

We can see that all the success rate for low weighted payload is higher than heavy weighted payload

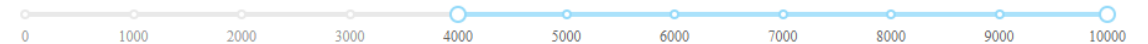
Payload range (Kg):



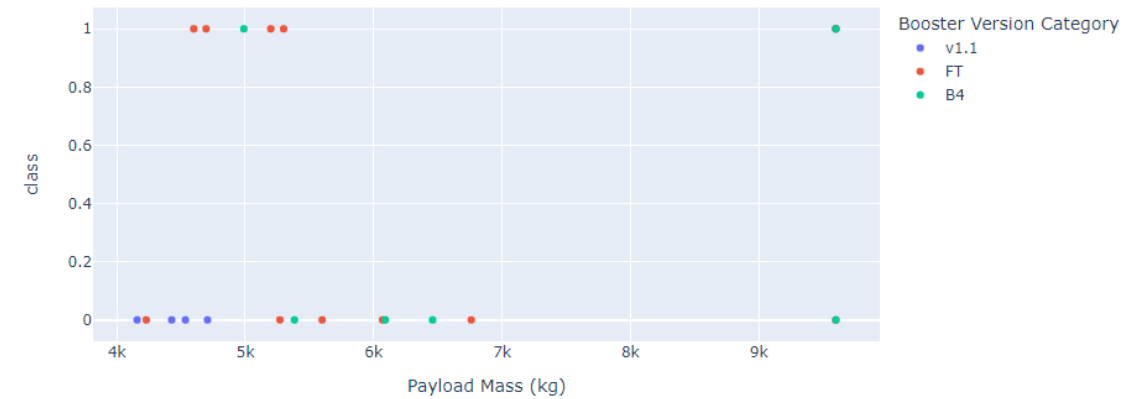
Correlation Between Payload and Success for All Sites



Payload range (Kg):



Correlation Between Payload and Success for All Sites



A magnifying glass is positioned over a bar chart. The chart has a light blue background and features two series of bars: blue and green. The x-axis is labeled with quarters: Q1, Q2, Q3, Q4. The y-axis has a label '1,000'. The magnifying glass is centered over the Q2 and Q3 bars, with the text 'Predictive Analysis' overlaid in white. The text 'Section 5' is also visible below the main title.

Predictive Analysis

Section 5

```
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 0.8892857142857145

Best Params is : {'criterion': 'entropy', 'max_depth': 6, 'max_features': 'sqrt', 'min_samples_leaf': 2, 'min_samples_split': 5, 'splitter': 'random'}

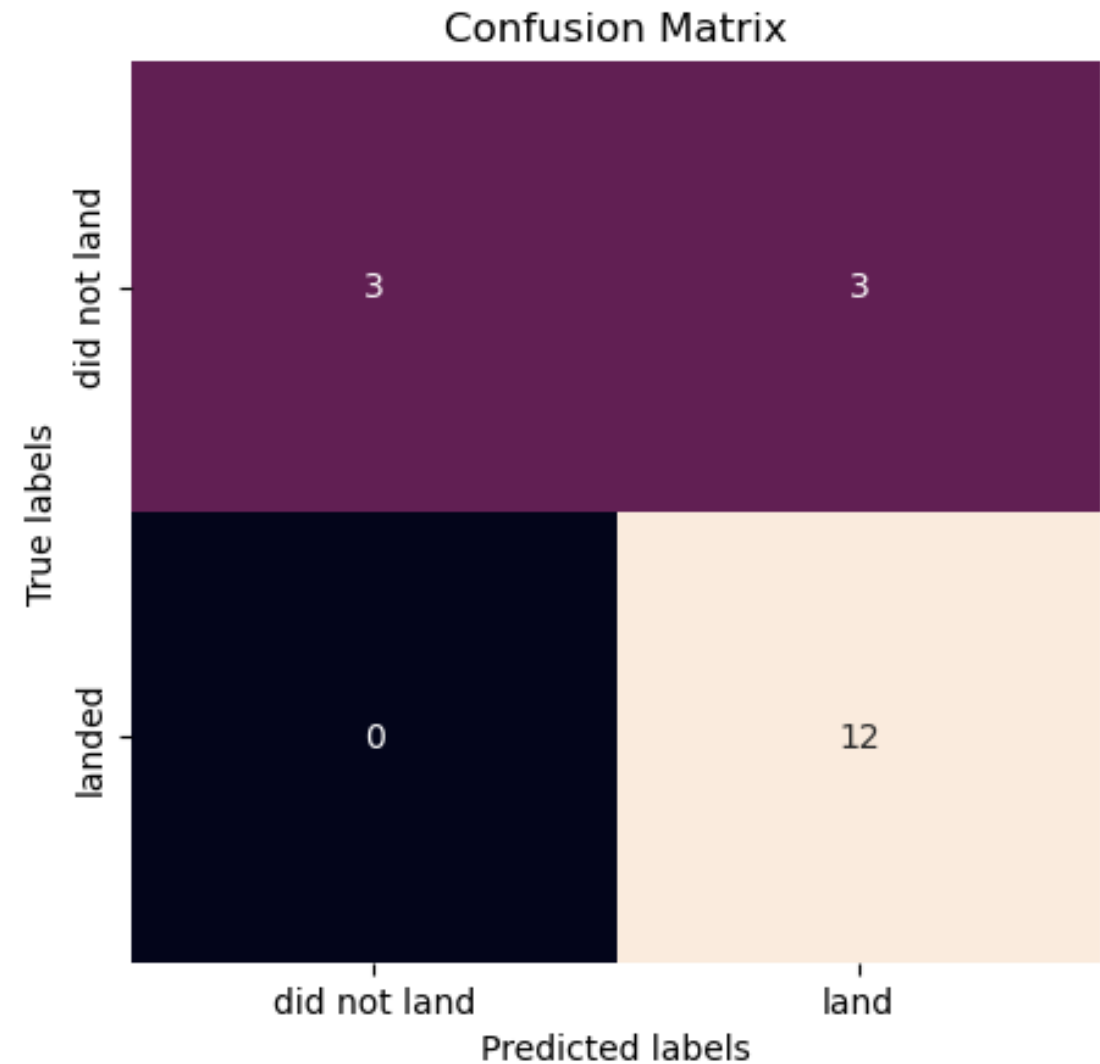
Classification Accuracy

- As we can see, by using the code as below: we could identify that the best algorithm to be the Tree Algorithm which have the highest classification accuracy.

Confusion Matrix

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

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



Conclusion

- The Tree Classifier Algorithm is the best Machine Learning approach 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



Thank You!

Thanks to:

[Instructors](#)

[Coursera](#)

[IBM](#)