



APPLIED DATA SCIENCE CAPSTONE

SPACEX

PROJECT

17TH JUNE 2024

GABRIEL GIDEON SIMDIMA

Outline

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



Executive Summary

Summary of methodologies

- Data collection
- Data wrangling
- EDA with data visualization
- EDA with SQL
- Building an interactive map with Folium
- Building a Dashboard with Plotly Dash
- Predictive analysis (Classification)

Summary of all

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

Introduction

SpaceX's achievements in rocket technology and launch operations have positioned them as a leader in the aerospace industry, with profound implications for the future of space exploration and commercial space travel through their robust Launch Sites and launch frequency.

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. If determining the landing of the first stage, the cost of a launch could be determined. The information obtained can be used by bidding alternate companies to compete against SpaceX for a rocket launch.

Consequently this project seek to identify the causal factors and relationship to cost and successful landing.

SECTION 1

Methodology



Methodology *Executive Summary*

Data collection from SpaceX Rest API, Web Scraping from page of Falcon and Falcon Heavy Langes

Data Wrangling performance through One Hot Encoding and removing empty dataset

Performing exploratory data analysis EDA and visualization using SQL

Performing interactive visual analytics using Folium and Plotly Dash

Perform predictive analysis using classification methods such as logistic regression, K-Nearest Neighbors, Support Vector Machines and Decision Tree

Data collection

In this project data collection was carried out through request SpaceX API and web scraping

1. The Request to SpaceX API consist of:

Gathering past data of SpaceX through source API after which it is retrieved, processed and ensuring that the data is the Falcon 9 launch. Also the missing payloads weight data from secret mission with average available

2. Web scraping

was carried out requesting Falcon 9 and Falcon Heavy launch data from Wikipedia page. Then accessed the Falcon 9 page, extracted all the column names from HTML Tables, parsed and transform into pandas frame ready for the analysis.

Data collection – spaceX API

- ❖ SpaceX API to collect data
- ❖ clean the requested data and data wrangling and formatting.
- ❖ Return spacex data to JSON
- ❖ Data into .csv

After response from API conver to a.JSON file

1

```
In [2]: 1 # Takes the dataset and uses the rocket column to call the API and append the data to the list
2 def getBoosterVersion(data):
3     for x in data['rocket']:
4         if x:
5             response = requests.get("https://api.spacexdata.com/v4/rockets/"+str(x)).json()
6             BoosterVersion.append(response['name'])
```

Custom function to clean data

2

```
[15]: 1 BoosterVersion
[15]: []

Now, let's apply getBoosterVersion function method to get the booster version

[16]: 1 # Call getBoosterVersion
2 getBoosterVersion(data)

the list has now been update

In [17]: 1 BoosterVersion[0:5]
Out[17]: ['Falcon 1', 'Falcon 1', 'Falcon 1', 'Falcon 1', 'Falcon 9']

In [18]: 1 # Call getLaunchSite
2 getLaunchSite(data)

In [19]: 1 # Call getPayloadData
2 getPayloadData(data)

In [20]: 1 # Call getCoreData
2 getCoreData(data)
```

List to Dict to dataframe

3

```
[21]: 1 launch_dict = {'FlightNumber': list(data['flight_number']),
2 'Date': list(data['date']),
3 'BoosterVersion':BoosterVersion,
4 'PayloadMass':PayloadMass,
5 'Orbit':Orbit,
6 'LaunchSite':LaunchSite,
7 'Outcome':Outcome,
8 'Flights':Flights,
9 'GridFins':GridFins,
10 'Reused':Reused,
11 'Legs':Legs,
12 'LandingPad':LandingPad,
13 'Block':Block,
14 'ReusedCount':ReusedCount,
15 'Serial':Serial,
16 'Longitude': Longitude,
17 'Latitude': Latitude}
18
```

Filter dataframe and export

4

```
[22]: 1 # Hint data['BoosterVersion']!='Falcon 1'
2 filt = df['BoosterVersion']!='Falcon 1'
3 data_falcon9 = df.loc[filt]
4 data_falcon9.head()
```

Out[24]:

Data Collection - Scraping

- ❖ We webscrap from Falcon 9 launch records using BeautifulSoup
- ❖ Then parsed the table and converted it into a pandas dataframe.
- ❖ Save dataframe to .csv

1

Creating BeautifulSoup object

```
In [7]: 1 # Use BeautifulSoup() to create a BeautifulSoup object from a response text content
2 soup = BeautifulSoup(response.content, 'html.parser')
```

2

Creating BeautifulSoup object

```
In [5]: 1 # use requests.get() method with the provided url
2 # assign the response to a object
3 response = requests.get(static_url)
```

3

Finding Table

```
In [8]: 1 # Use the find_all function in the BeautifulSoup object, with element name
2 # Assign the result to a list called 'html_tables'
3 html_tables = soup.find_all('table')
```

4

Getting Columns Name

```
In [10]: 1 column_names = []
2
3 # Apply find_all() function with 'th' element
4 # Iterate each th element and apply the provided function
5 # Append the Non-empty column name ('if name is not None')
6
7 table = first_launch_table.find_all('th')
8 for row in table:
9     name = extract_column_from_header(row)
10    if name is not None and len(name) > 0:
11        column_names.append(name)
```

5

Getting Columns Name

```
launch_dict = dict.fromkeys(column_names)
2
3 # Remove an irrelevant column
4 del launch_dict['Date and time ( )']
5
6 # Let's initial the launch_dict with each column name
7 launch_dict['Flight No.'] = []
8 launch_dict['Launch site'] = []
9 launch_dict['Payload'] = []
10 launch_dict['Payload mass'] = []
11 launch_dict['Orbit'] = []
12 launch_dict['Customer'] = []
13 launch_dict['Launch outcome'] = []
14 # Added some new columns
15 launch_dict['Version Booster'] = []
16 launch_dict['Booster landing'] = []
17 launch_dict['Date'] = []
18 launch_dict['Time'] = []
```

6

Append data to keys

```
[13]: 1 extracted_row = 0
2 #Extract each table
3 for table_number, table in enumerate(html_tables):
4     # get table row
5     for rows in table.find_all("tr"):
```

7

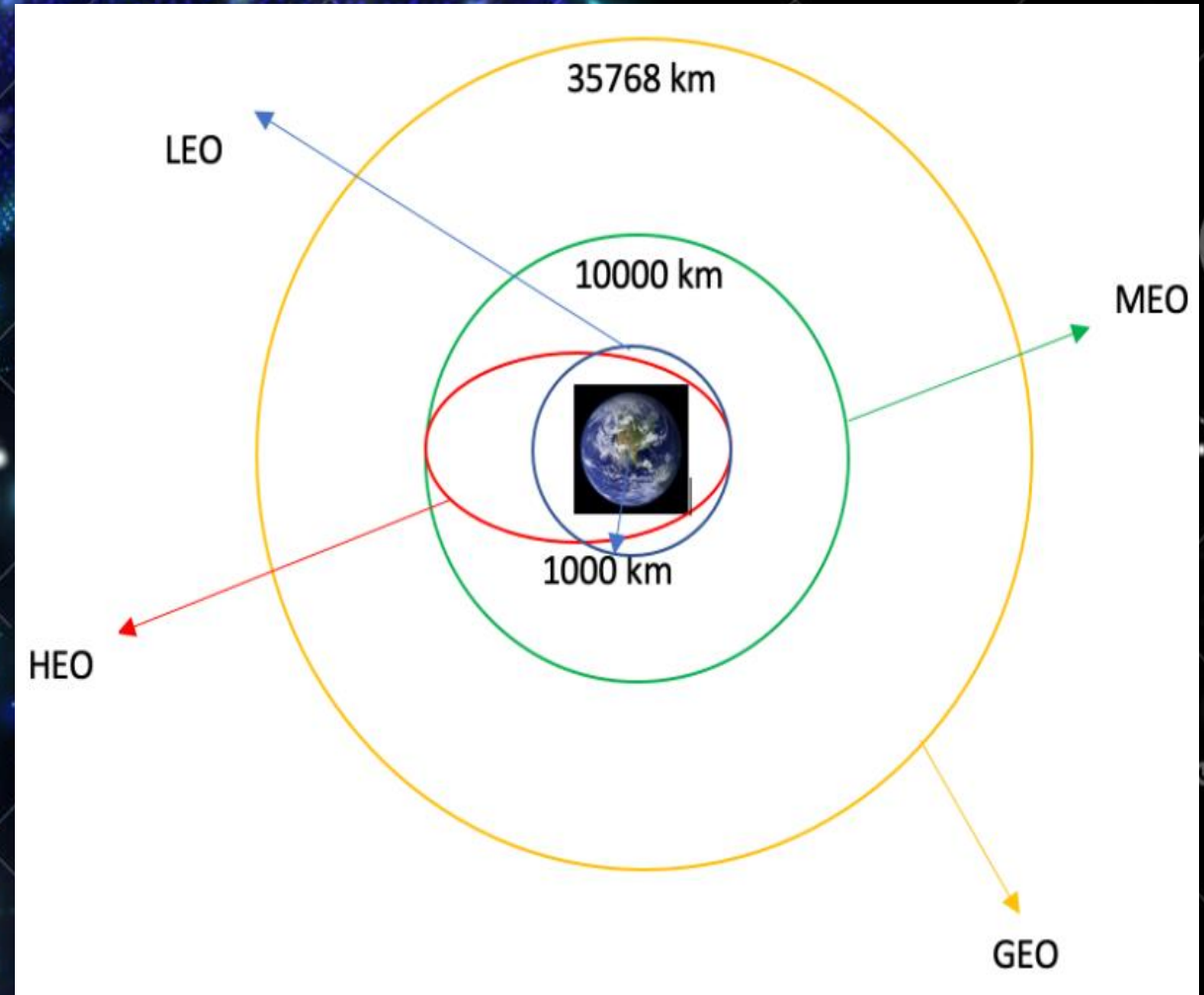
Covert dict to dataframe and to .csv

```
1 df = pd.DataFrame(launch_dict)
2 df.head()

[16]: 1 df.to_csv('spacex_web_scraped.csv', index=False)
```


Data Wrangling

- ❖ The number of launches at each site and number of occurrences of each orbit were determined
- ❖ Landing outcome label frame outcome table and the result in .csv was exported.



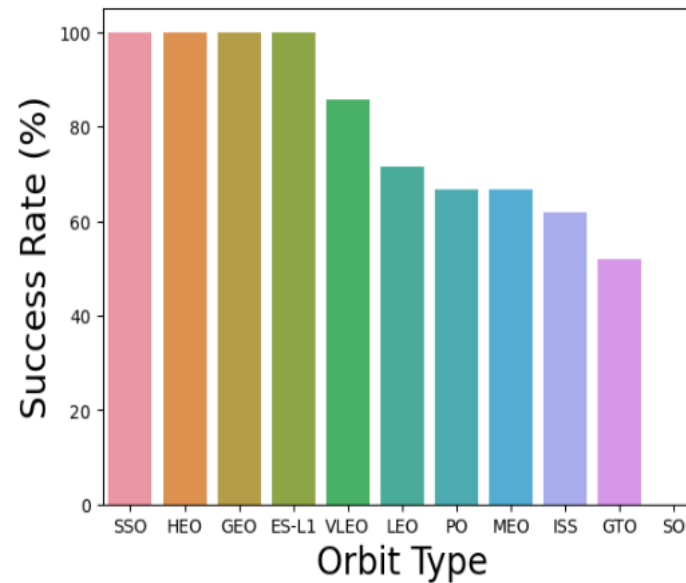
Eda with Data Visualization

Visualized the relationship between

- ❖ flight number and Payload mass,
- ❖ Payload and launch site
- ❖ Orbit and flight number
- ❖ Flight number and Launch Site
- ❖ Orbit and Payload Mass
- ❖ Payload and Orbit Type

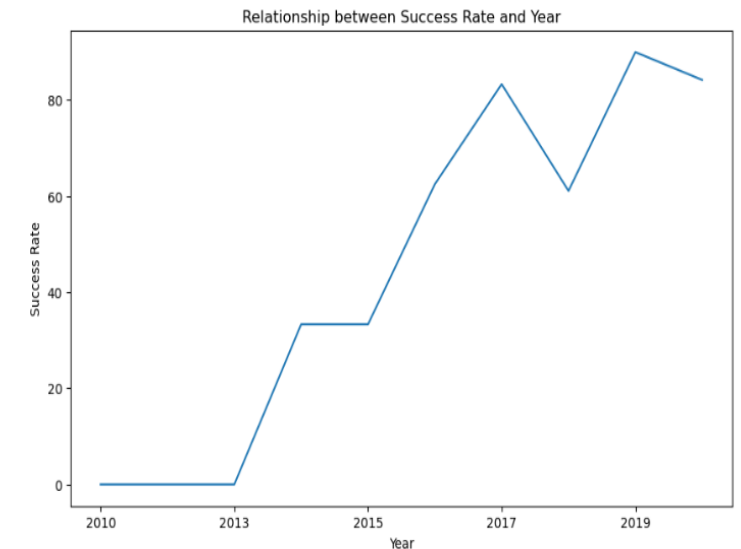
Mean vs Orbit

```
[39]: 1 # HINT use groupby method on Orbit column and get the mean of Class column
      2 orbit_success = df.groupby(['Orbit'])['Class'].aggregate(np.average).reset
      3 orbit_success['Class'] = np.round(orbit_success['Class']*100, 2)
      4 sns.barplot(x='Orbit', y='Class', data=orbit_success)
      5 plt.ylabel("Success Rate (%)", fontsize=20)
      6 plt.xlabel("Orbit Type", fontsize=20)
      7 plt.show()
```



Success rate vs Year

```
[16]: 1 # Plotting a Line chart with x axis to be the extracted year and y axis to be the success rate
      2 Extract_year(df)
      3 df['Year']=year
      4 fig,ax=plt.subplots()
      5 df_success1=df.groupby('Year')['Class'].mean()*100
      6 df_success1.plot(kind='line', figsize=(10,6))
      7 plt.xlabel('Year') # add to x-label to the plot
      8 plt.ylabel('Success Rate') # add y-label to the plot
      9 plt.title('Relationship between Success Rate and Year') # add title to the plot
      10
      11 plt.show()
```



EDA with SQL

1 SpaceX dataset

We loaded the SpaceX dataset into a PostgreSQL database without leaving the jupyter notebook.

2 We applied EDA

EDA with SQL were applied to get insight from the data. We wrote queries to find out:

- ❖ The names of unique launch sites in the space mission.
- ❖ The total payload mass carried by boosters launched by NASA (CRS)
- ❖ The average payload mass carried by booster version F9 v1.1
- ❖ The total number of successful and failure mission outcomes
- ❖ The failed landing outcomes in drone ship, their booster version and launch site names

Interactive Map with Folium

The activities carried out are as follow

- ❖ 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.

Dashboard with Plotly Dash

The activities carried out are

- ❖ We built an interactive dashboard with Plotly dash
- ❖ We plotted pie charts showing the total launches by a certain sites
- ❖ We plotted scatter graph showing the relationship with Outcome and Payload Mass (Kg) for the different booster version.

Predictive Analysis (Classification)

Here we

- ❖ loaded the data using numpy and pandas, then transformed the data, split the data into training and testing.
- ❖ built different machine learning models and tune different hyperparameters using GridSearchCV.
- ❖ used accuracy as the metric for out the model, improved the model using feature engineering and algorithm tuning.
- ❖ lastly determined the best performing classification model.

Results

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

SECTION 2

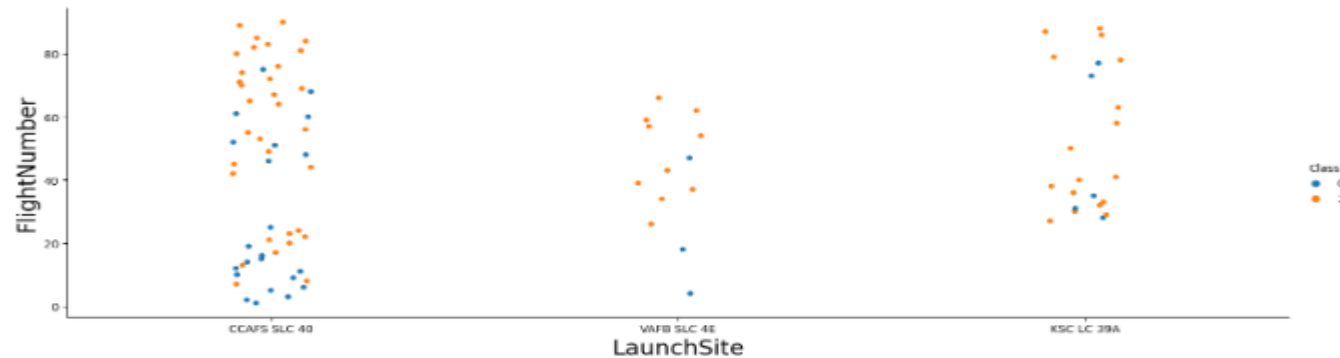
Insight Drawn from EDA



Flight Number Vs Lunch Site

```
1 # Plot a scatter point chart with x axis to be Flight Number and y axis to be
2 sns.catplot(y="FlightNumber",x="LaunchSite",hue='Class',data=df, aspect=3)
3 plt.xlabel("LaunchSite",fontsize=20)
4 plt.ylabel("FlightNumber",fontsize=20)
5 plt.show()
```

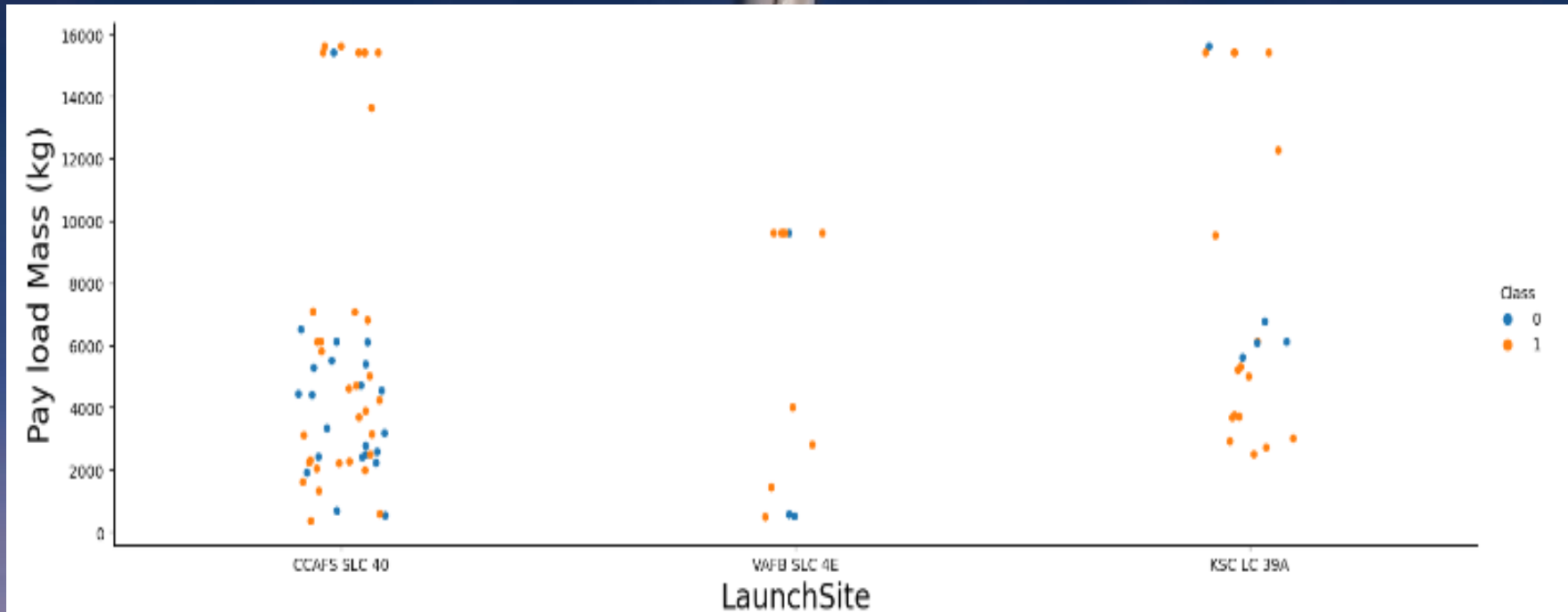
C:\Users\GABRIEL SIMDIMA\anaconda3\Lib\site-packages\seaborn\axisgrid.py:118: UserWarning: The figure layout has changed to tight
self._figure.tight_layout(*args, **kwargs)



Findings:

1. CCAFS LC-40: No strong relationship between flight number and success rate. Success appears consistent across various flight numbers.
2. KSC LC-39A: Similarly, no strong correlation between flight number and success rate. Success rates vary across flight numbers.
3. VAFB SLC 4E: Success rates appear consistent, regardless of flight number.

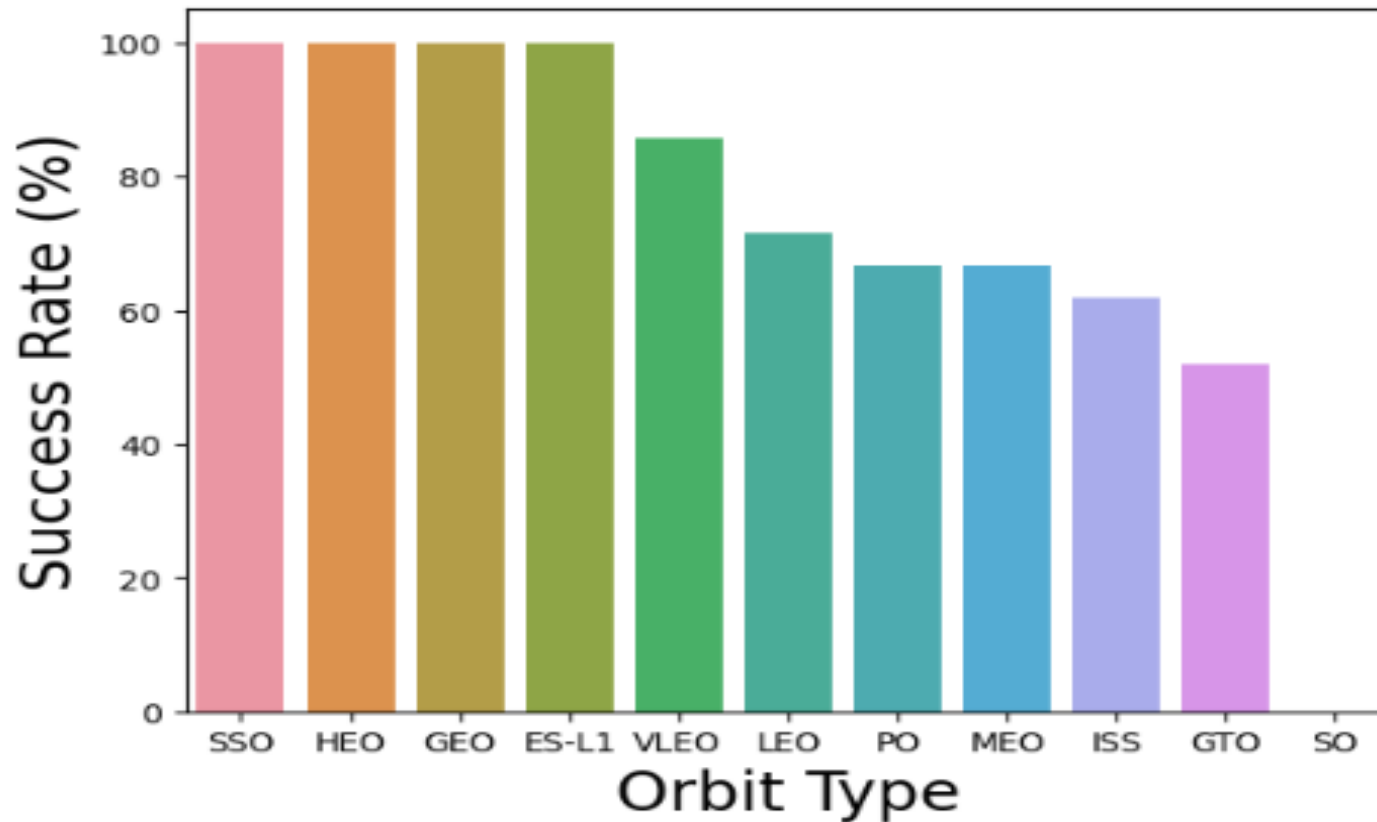
Payload Mass Vs Launch Site



Findings:

1. VAFB-SLC: No rockets launched for heavy payload mass (greater than 10,000).
2. Other Launch Sites: Rockets launched for a range of payload masses.
3. The launch site affects payload capacity and handling.

Success Rate Vs Orbit Type



Findings:

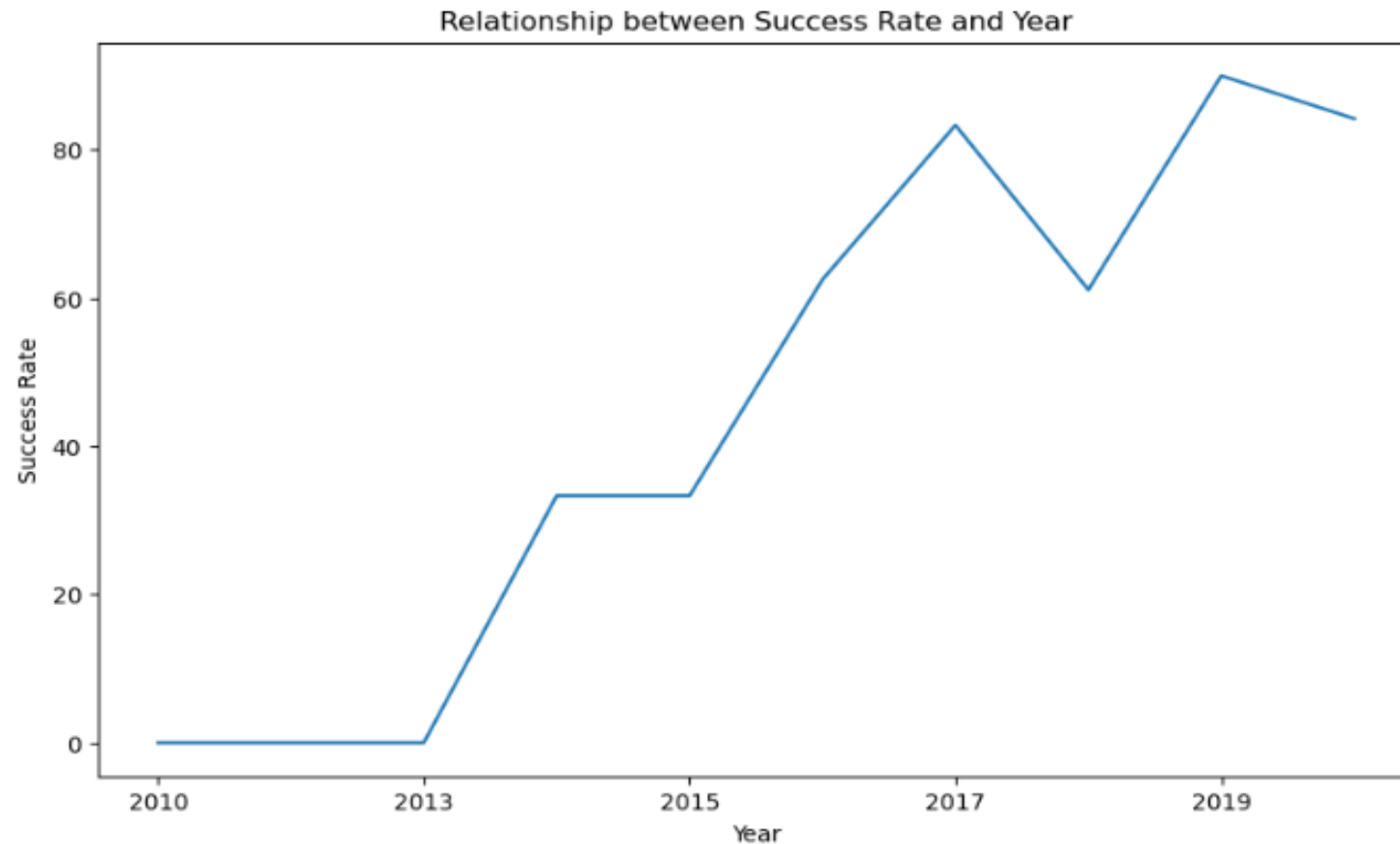
1. SSE, HEO, GEO, ES-L1 Orbits: 100% success rates, indicating high reliability.
2. Other Orbit Types: Varying success rates, suggesting mission complexities.

Findings:

1. LEO Orbit: Success appears correlated with the number of flights.
2. GTO Orbit: No clear relationship between flight number and success.

1. PO, LEO, ISS Orbits: Higher payload mass corresponds to a more successful landing.
2. GTO Orbit: Payload mass does not clearly affect landing success.

Success Rate Vs Year



Findings:

1. 2013 to 2020: There is consistent increase in success rates.
2. Highlight: SpaceX's continual improvement in launch reliability.

Unique Lunch Sites

This shows
the unique
Lunch Sites
from the
spaceX data
obtained

```
In [11]: 1 count=0
          2 for site in df["Launch_Site"]:
          3     if "CCA" in site and count < 5:
          4         print(site)
          5         count=count+1
          6
```

```
CCAFS LC-40
CCAFS LC-40
CCAFS LC-40
CCAFS LC-40
CCAFS LC-40
```

SECTION 3

EDA with .SQL



Unique Lunch Sites

This shows
the unique
Lunch Sites
from the
spaceX data
obtained

```
In [11]: 1 count=0
          2 for site in df["Launch_Site"]:
          3     if "CCA" in site and count < 5:
          4         print(site)
          5         count=count+1
          6
```

```
CCAFS LC-40
CCAFS LC-40
CCAFS LC-40
CCAFS LC-40
CCAFS LC-40
```

Lunch Sites names beginning with CCA

```
In [48]: 1 %sql SELECT * \
        2 FROM SPACEXTBL \
        3 WHERE LAUNCH_SITE LIKE 'CCA%' LIMIT 5;
```

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

```
Done.
```

```
Out[48]:
```

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

Total payload
carried by boosters
from NASA as
45596 AS shown

```
In [49]: 1 %sql SELECT SUM(PAYLOAD_MASS_KG_) \
          2 FROM SPACEXTBL \
          3 WHERE CUSTOMER = 'NASA (CRS)';
```

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

```
Done.
```

```
Out[49]: SUM(PAYLOAD_MASS_KG_)
```

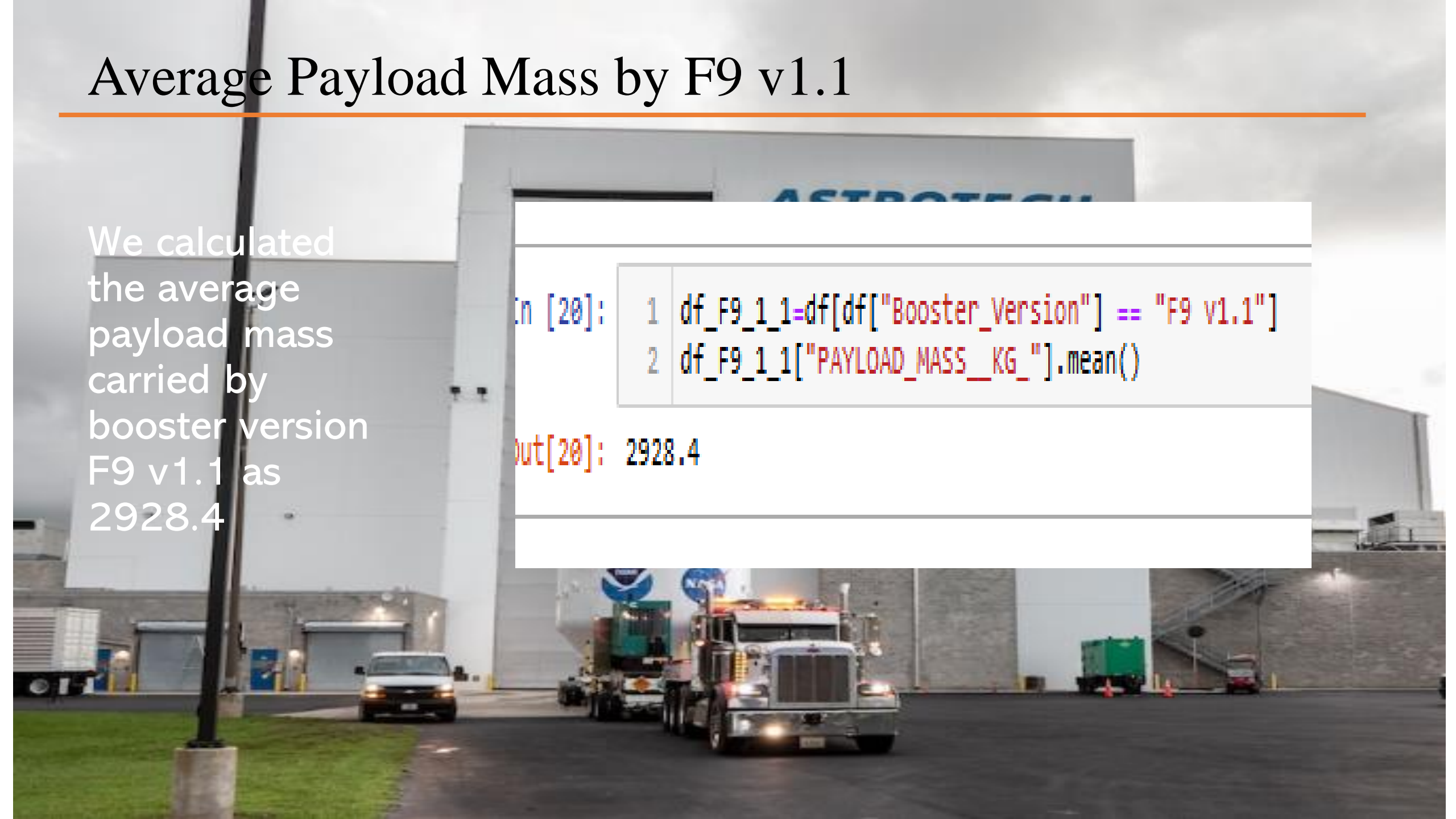
45596

Average Payload Mass by F9 v1.1

We calculated the average payload mass carried by booster version F9 v1.1 as 2928.4

```
In [20]: 1 df_F9_1_1=df[df["Booster_Version"] == "F9 v1.1"]  
        2 df_F9_1_1["PAYLOAD_MASS_KG_"].mean()
```

```
Out[20]: 2928.4
```



First successful landing outcome Date

```
In [21]: 1 df_ground_pad=df[df['Landing_Outcome'] == 'Success (ground pad)']  
        2 df_ground_pad.iloc[0,0]
```

```
Out[21]: '22/12/2015'
```

Successful Drone Ship Landing with Payload between 4000 and 6000

```
In [22]: 1 %sql select BOOSTER_VERSION from SPACEXTBL where Landing_Outcome='Success (drone ship)' and PAYLOAD_MASS_KG BETWEEN 4000 a
```

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

```
Out[22]:
```

Booster_Version
F9 FT B1022
F9 FT B1028
F9 FT B1021.2
F9 FT B1031.2

Total Number of Successful and Failure Mission Outcomes

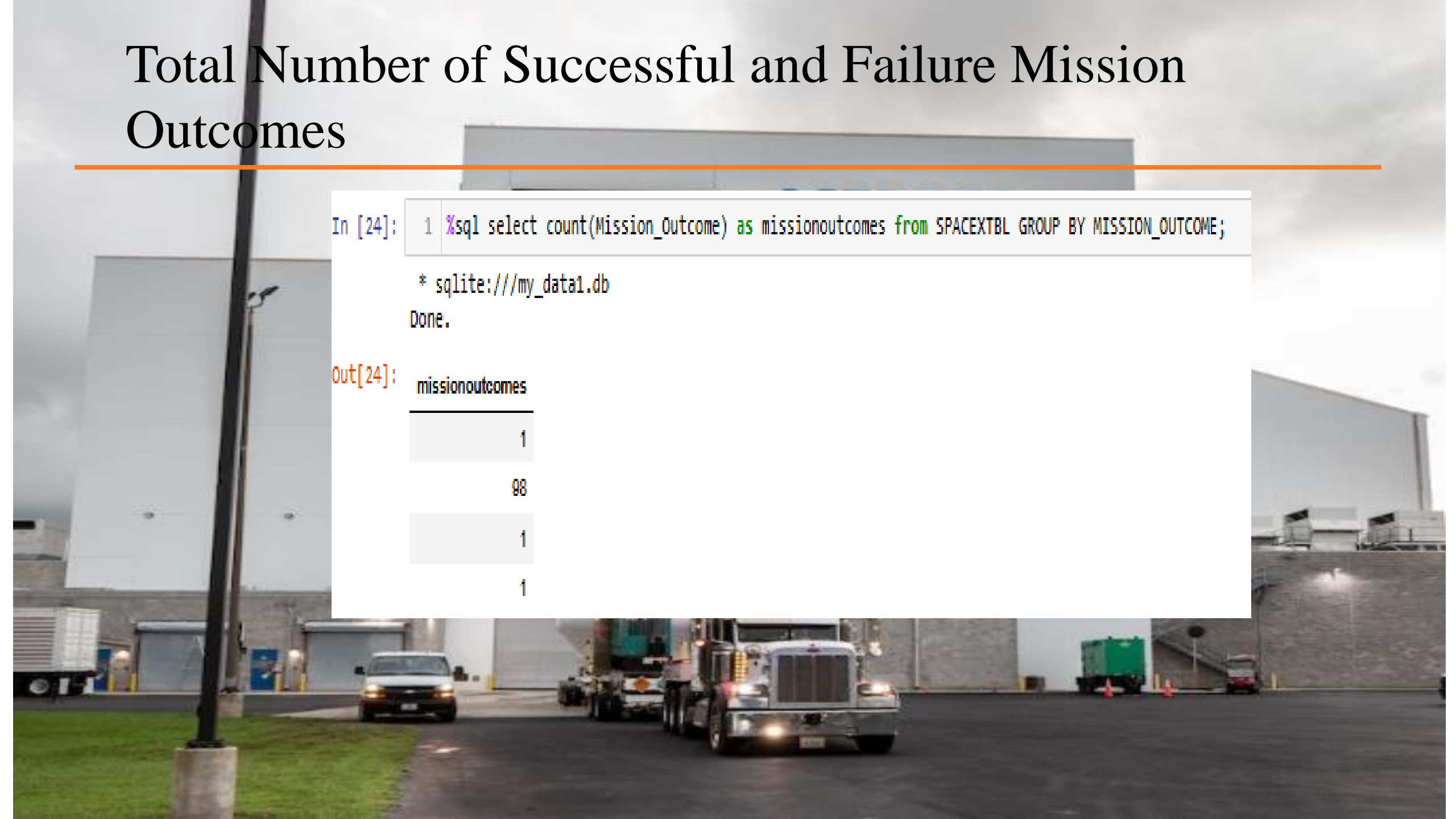
```
In [24]: 1 %sql select count(Mission_Outcome) as missionoutcomes from SPACEXTBL GROUP BY MISSION_OUTCOME;
```

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

```
Done.
```

```
Out[24]: missionoutcomes
```

1
98
1
1



Boosters Carried Maximum Payload

```
In [26]: 1 %sql select BOOSTER_VERSION as boosterversion from SPACEXTBL where PAYLOAD_MASS_KG_=(select max(PAYLOAD_MASS_KG_) from
```

* sqlite:///my_data1.db
Done.

```
Out[26]:
```

boosterversion
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

```
In [29]: 1 %sql SELECT substr(Date,4,2) as month, Date,Booster_Version, Launch_Site, [Landing_Outcome] \
2 FROM SPACEXTBL \
3 where [Landing_Outcome] = 'Failure (drone ship)' and substr(Date,7,4)='2015';
```

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

Done.

```
Out[29]:
```

	month	Date	Booster_Version	Launch_Site	Landing_Outcome
	01	10/01/2015	F9 v1.1 B1012	CCAFS LC-40	Failure (drone ship)
	04	14/04/2015	F9 v1.1 B1015	CCAFS LC-40	Failure (drone ship)

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

```
[31]: 1 %sql SELECT [Landing_Outcome], count(*) as count_outcomes \
      2 FROM SPACEXTBL \
      3 WHERE DATE between '04-06-2010' and '20-03-2017' group by [Landing_Outcome] order by count_outcomes DESC
```

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

```
[31]:
```

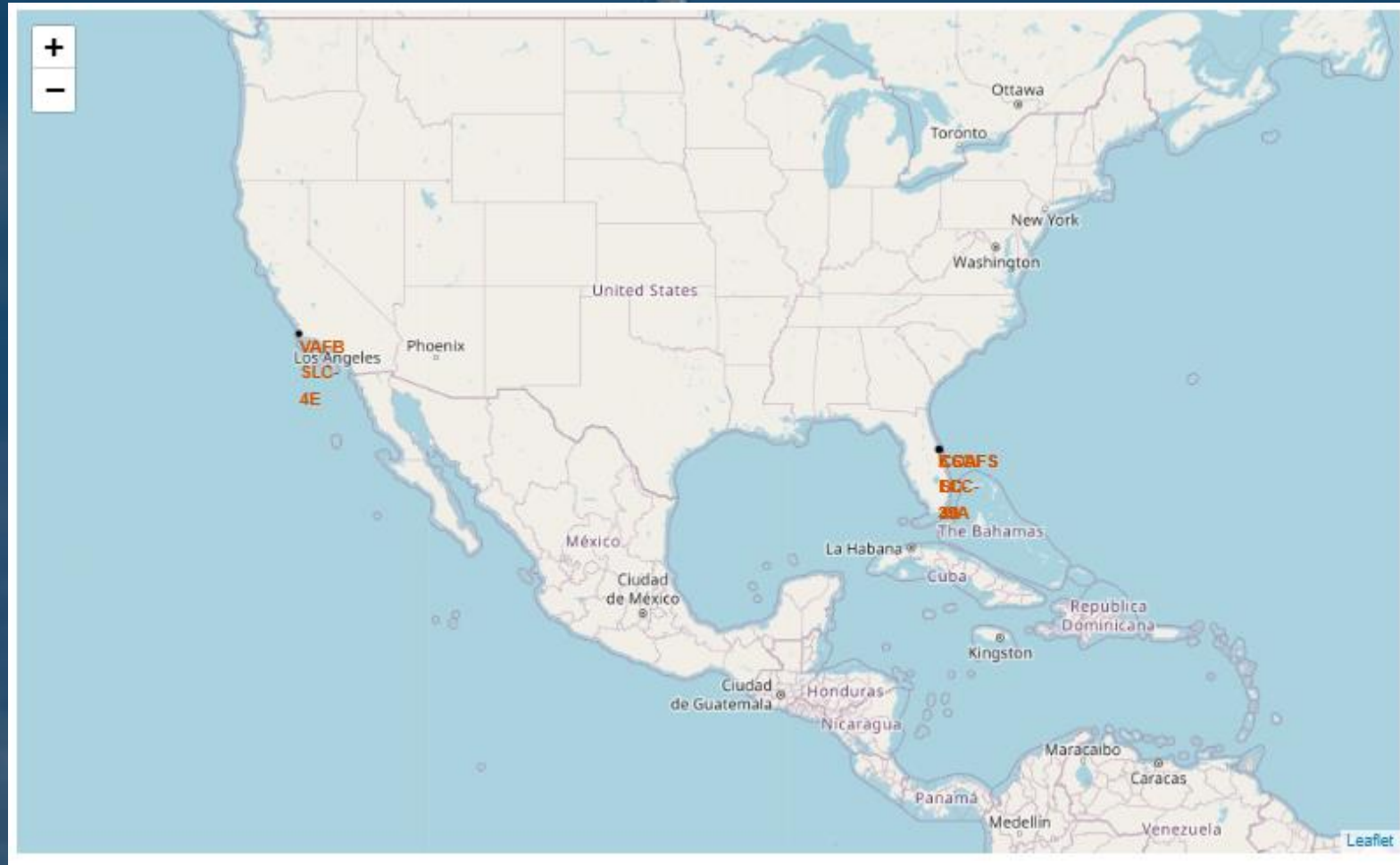
Landing_Outcome	count_outcomes
Success	21
No attempt	10
Success (drone ship)	8
Success (ground pad)	6
Failure (drone ship)	5
Failure	3
Controlled (ocean)	3
Failure (parachute)	2
No attempt	1

SECTION 4

Launch Sites Proximities Analysis



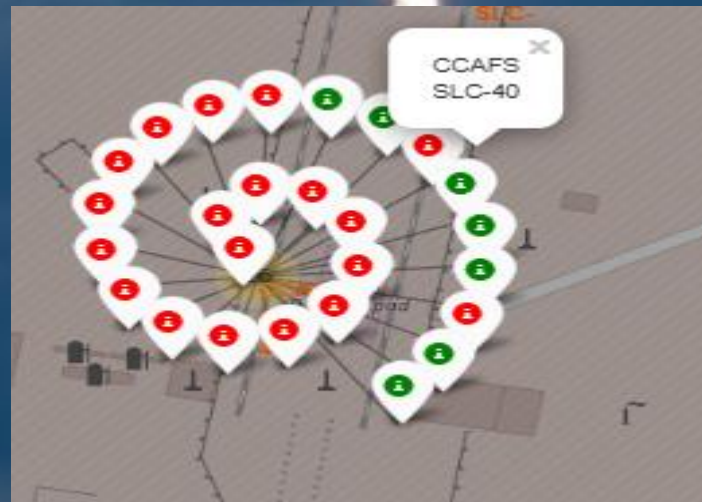
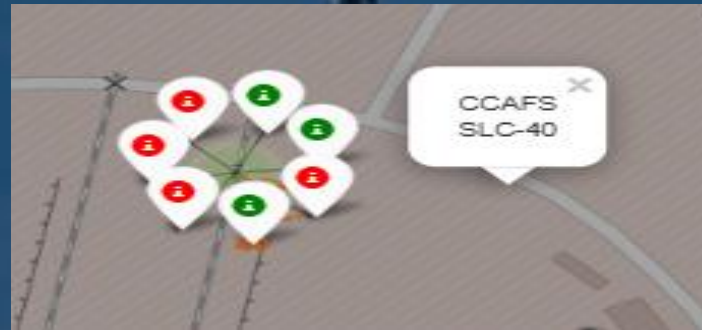
All launch sites global map markers



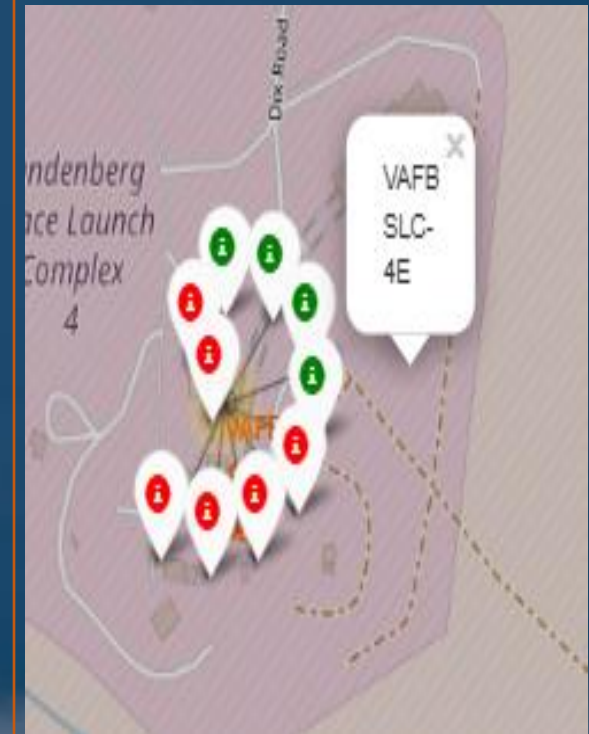
All launch sites global map markers



Florida launch site with green indicating successful while red indicating unsuccessful launch

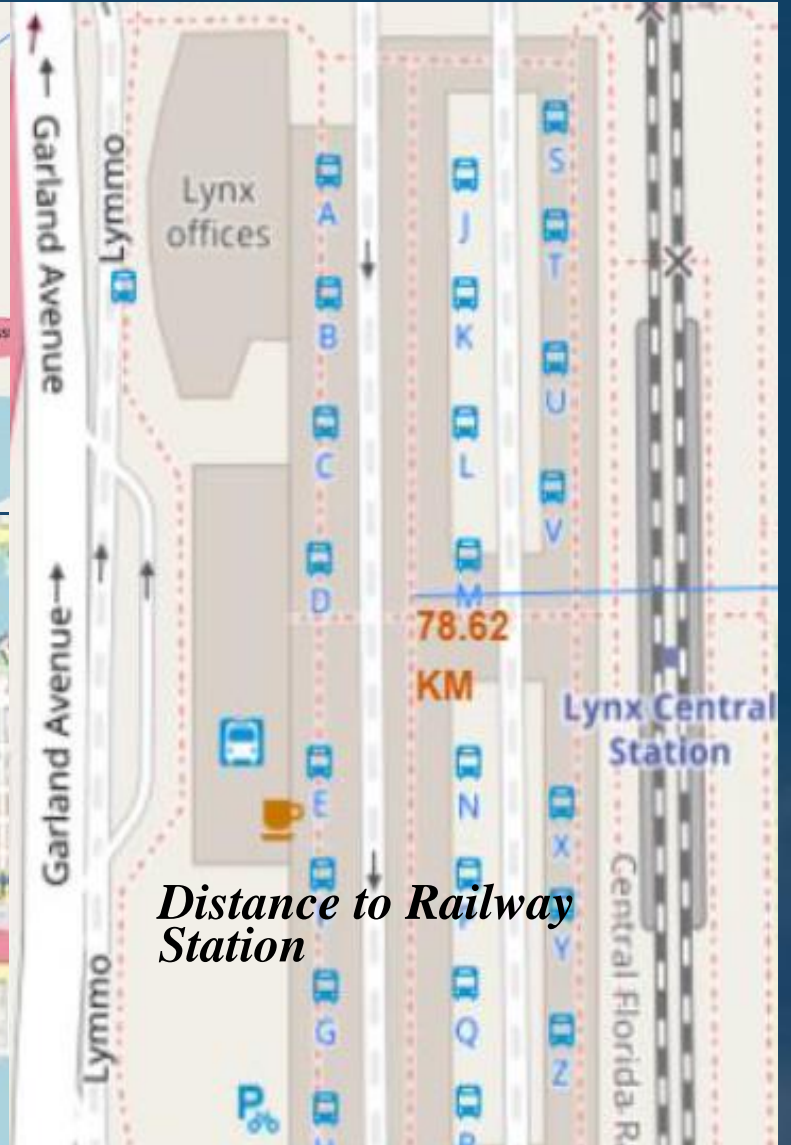
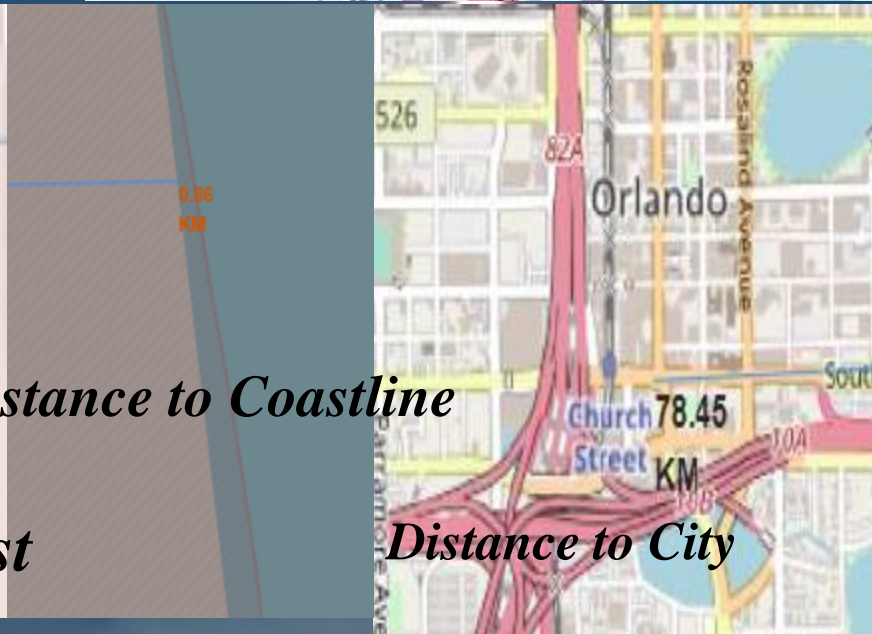


California launch site



Launch Site distance to landmarks

- Are the launch sites in close proximity to railway? NO
- Are the launch sites in close proximity to highway? NO
- Are the launch sites in close proximity to Coaster line? YES
- Do launch sites keep certain distance away from the city? YES



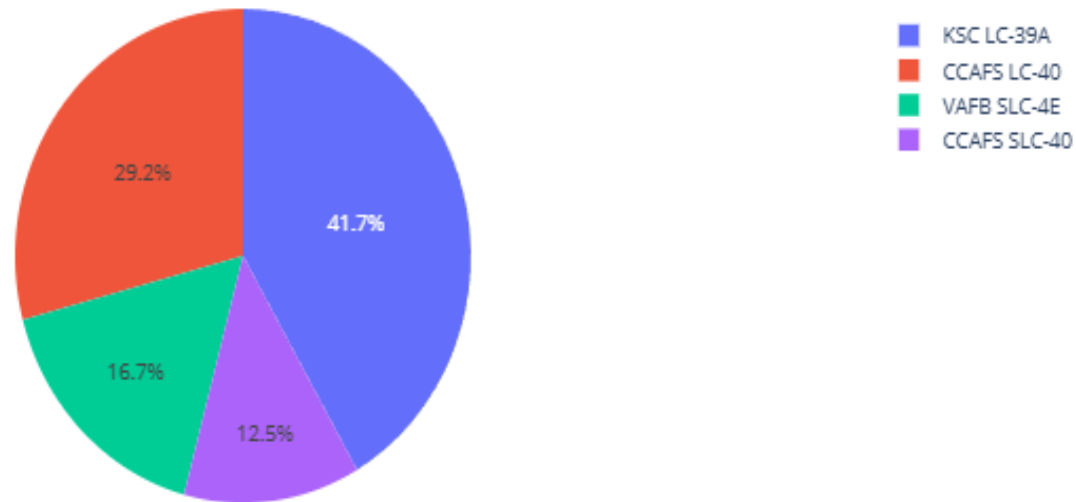
SECTION 5

Build a Dashboad with Plotly Dash

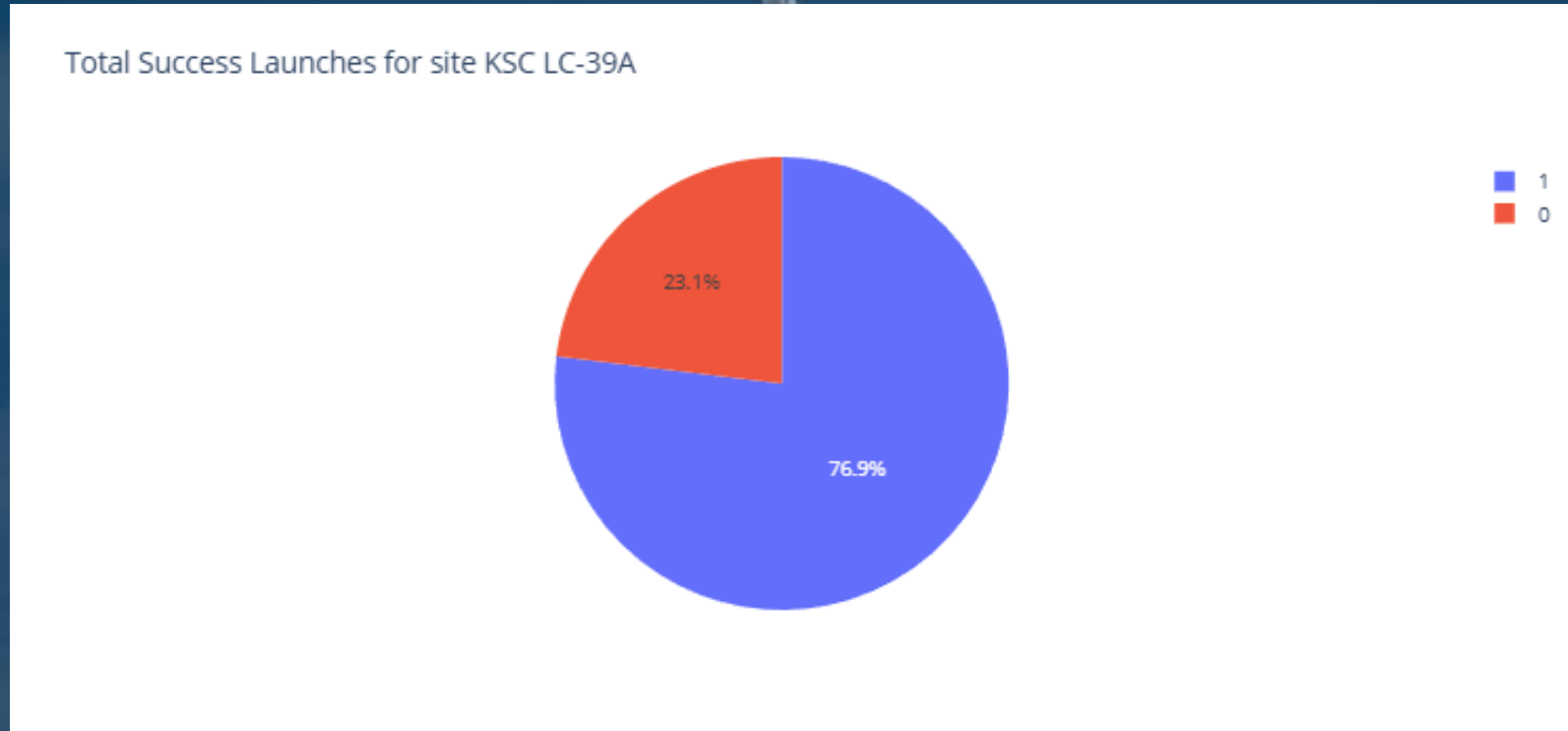


Pie chart showing the Launch site with the highest launch success ratio

Success Count for all launch sites

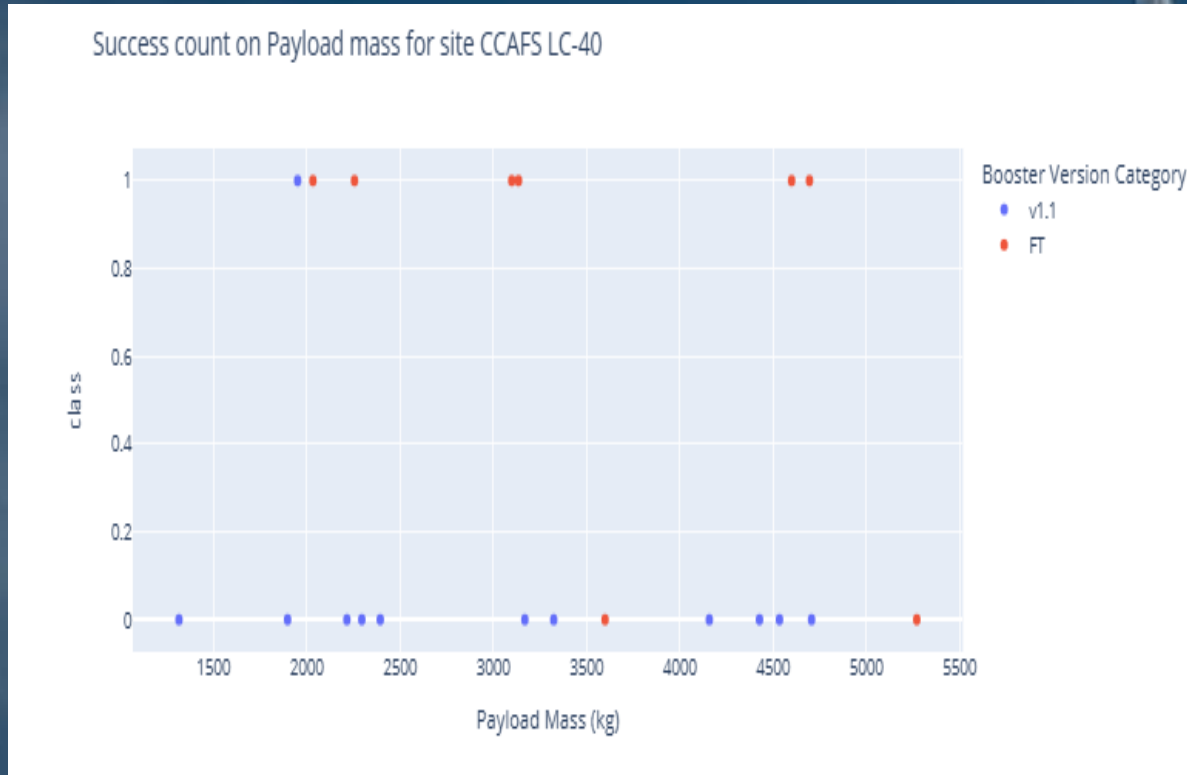


Pie chart showing the Launch site with the highest launch success ratio

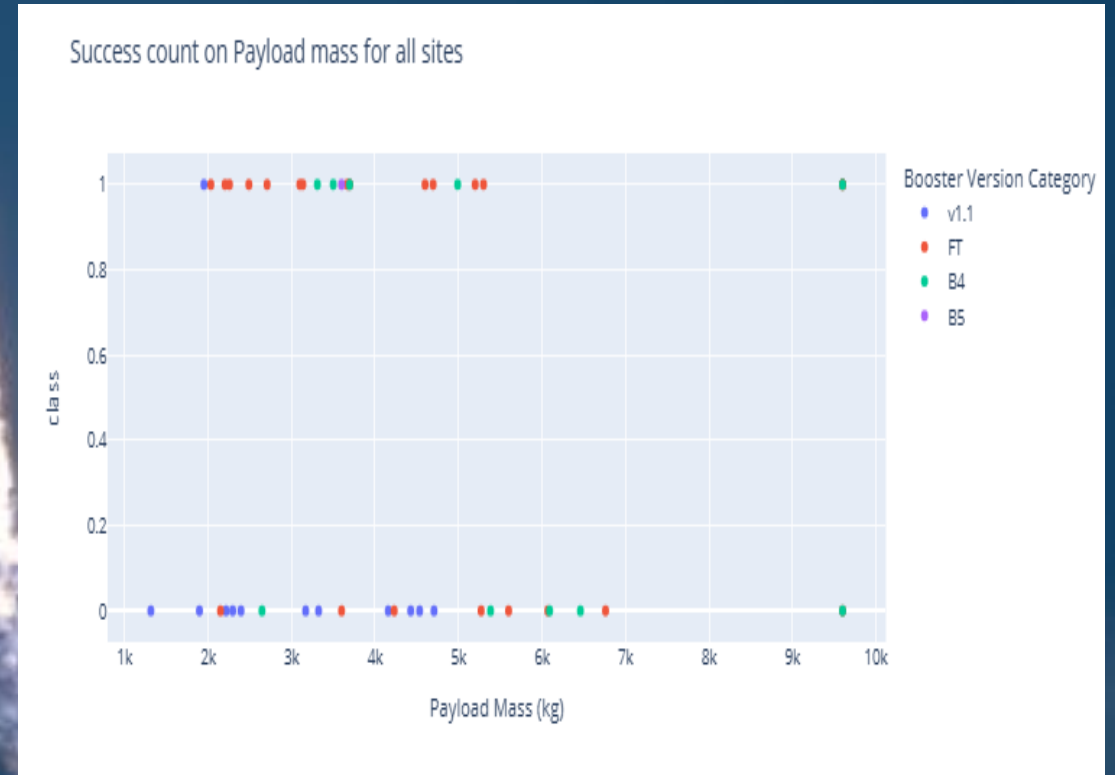


Scatter plot of Payload vs Launch Outcome for all sites, with different payload selected in the range slider

Low weighted payload 0-4000kg



Heavy weighted payload 4000-10000kg



We can also conclude therefore that the success of the low weighted payloads is higher than the heavy weighted payloads

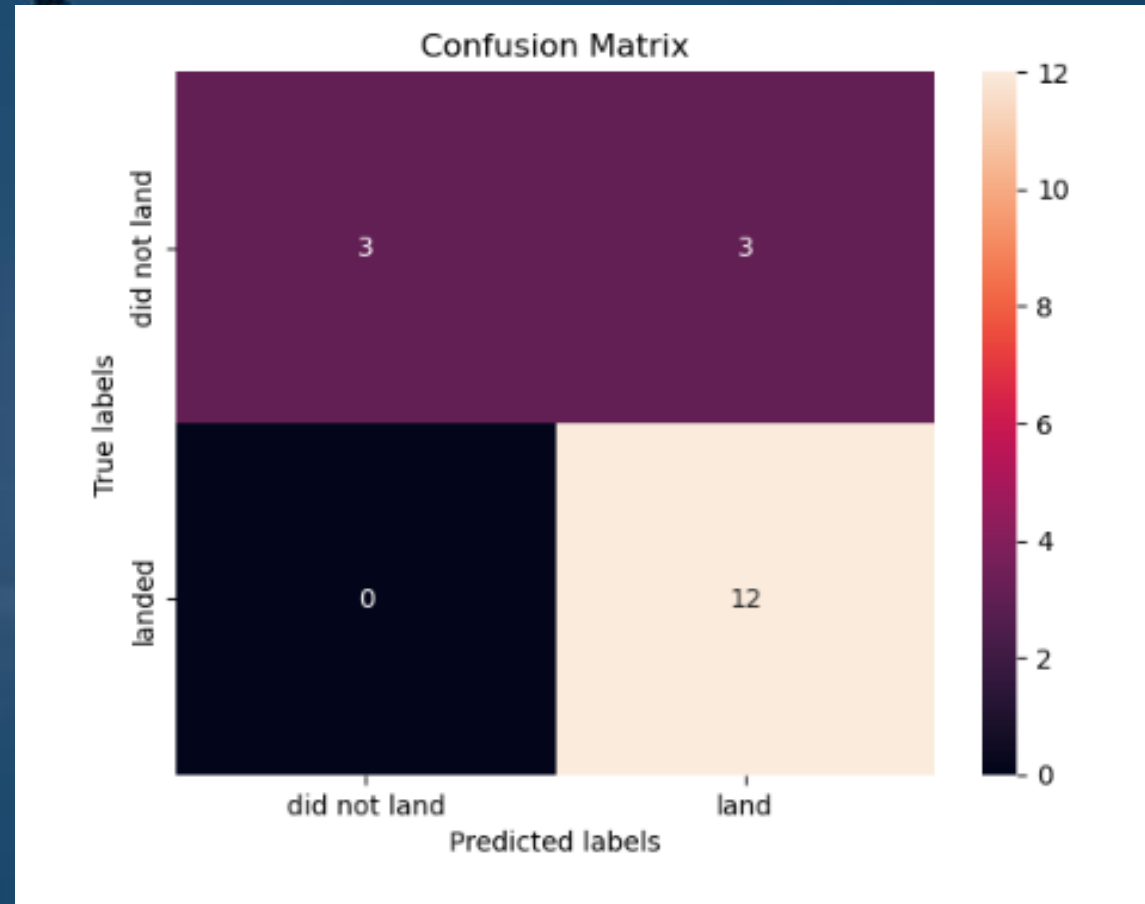
SECTION 6

Predictive Analysis (Classification)



Confusion Matrix

From confusion matrix results of the prediction using KNN, Logistic Regression, SVM and Decision Tree model, the major problem is the false positives .i.e., unsuccessful landing marked as successful landing by the classifier.



Classification Accuracy using Training Data

From the result shown the prediction accuracy

However Decision Tree is the least with Accuracy of 72.2%

```
1 Report = pd.DataFrame({'Method' : ['Test Data Accuracy']})
2
3 knn_accuracy=knn_cv.score(X_test, Y_test)
4 Decision_tree_accuracy=tree_cv.score(X_test, Y_test)
5 SVM_accuracy=svm_cv.score(X_test, Y_test)
6 Logistic_Regression=logreg_cv.score(X_test, Y_test)
7
8 Report['Logistic_Reg'] = [Logistic_Regression]
9 Report['SVM'] = [SVM_accuracy]
10 Report['Decision Tree'] = [Decision_tree_accuracy]
11 Report['KNN'] = [knn_accuracy]
12
13 Report.transpose()
```

0	
Method	Test Data Accuracy
Logistic_Reg	0.833333
SVM	0.833333
Decision Tree	0.722222
KNN	0.833333

Conclusion

By way of conclusion we can conclude that:

1. The larger the flight amount at a launch site, the greater the success rate at a launch site.
2. Launch success rate started to increase in 2013 till 2020.
3. The Orbits ES-L1, GEO, HEO, SSO, VLEO had the most success rate.
4. The KSC LC-39A had the most successful launches of any sites.
5. The Decision tree classifier seems to be the least of the machine learning algorithm for this task, though the result show it is a good method for this task.

Thank you very much.

