

Winning Space Race with Data Science

Ujjwala Kumari Singh
7th January 2023



Outline

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

Executive Summary

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

Introduction

- Project background and context

Space X 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 Space X 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 space X for a rocket launch. This goal of the project is to create a machine learning pipeline to predict if the first stage will land successfully.

- Problems you want to find answers

- What factors determine if the rocket will land successfully?
- The interaction amongst various features that determine the success rate of a successful landing.
- What operating conditions needs to be in place to ensure a successful landing program.

Section 1

Methodology

Methodology

Executive Summary

- Data collection methodology:
 - Data was collected using SpaceX API and web scraping from Wikipedia.
- Perform data wrangling
 - One-hot encoding was applied to categorical features
- 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

- The data was collected using various methods
 - Data collection was done using get request to the SpaceX API.
 - Next, we decoded the response content as a Json using `.json()` function call and turn it into a pandas dataframe using `.json_normalize()`.
 - We then cleaned the data, checked for missing values and fill in missing values where necessary.
 - In addition, we performed web scraping from Wikipedia for Falcon 9 launch records with BeautifulSoup.
 - The objective was to extract the launch records as HTML table, parse the table and convert it to a pandas dataframe for future analysis.

Data Collection – SpaceX API

- We used the get request to the SpaceX API to collect data, clean the requested data and did some basic data wrangling and formatting
- The link to the notebooks is “<https://github.com/ujjwalasingh22/IB-M-data-science/blob/master/Data%20wrangling.ipynb>”

1. Get request for rocket launch data using API

```
In [6]: spacex_url="https://api.spacexdata.com/v4/launches/past"
```

```
In [7]: response = requests.get(spacex_url)
```

2. Use json_normalize method to convert json result to dataframe

```
In [12]: # Use json_normalize method to convert the json result into a dataframe  
# decode response content as json  
static_json_df = res.json()
```

```
In [13]: # apply json_normalize  
data = pd.json_normalize(static_json_df)
```

3. We then performed data cleaning and filling in the missing values

```
In [30]: rows = data_falcon9['PayloadMass'].values.tolist()[0]  
  
df_rows = pd.DataFrame(rows)  
df_rows = df_rows.replace(np.nan, PayloadMass)  
  
data_falcon9['PayloadMass'][0] = df_rows.values  
data_falcon9
```

Data Collection - Scraping

- We applied web scrapping to webscrap Falcon 9 launch records with BeautifulSoup
- We parsed the table and converted it into a pandas dataframe.
- The link to the notebook is “<https://github.com/ujjwalasingh22/IBM-data-science/blob/master/Data%20wrangling.ipynb>”

```
1. Apply HTTP Get method to request the Falcon 9 rocket launch page
```

```
In [4]: static_url = "https://en.wikipedia.org/w/index.php?title=List_of_Falcon_9_and_Falcon_Heavy_launches&oldid=1027686922"
```

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

```
Out[5]: 200
```

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

```
In [6]: # Use BeautifulSoup() to create a BeautifulSoup object from a response text content  
soup = BeautifulSoup(html_data.text, 'html.parser')
```

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

```
In [7]: # Use soup.title attribute  
soup.title
```

```
Out[7]: <title>List of Falcon 9 and Falcon Heavy launches - Wikipedia</title>
```

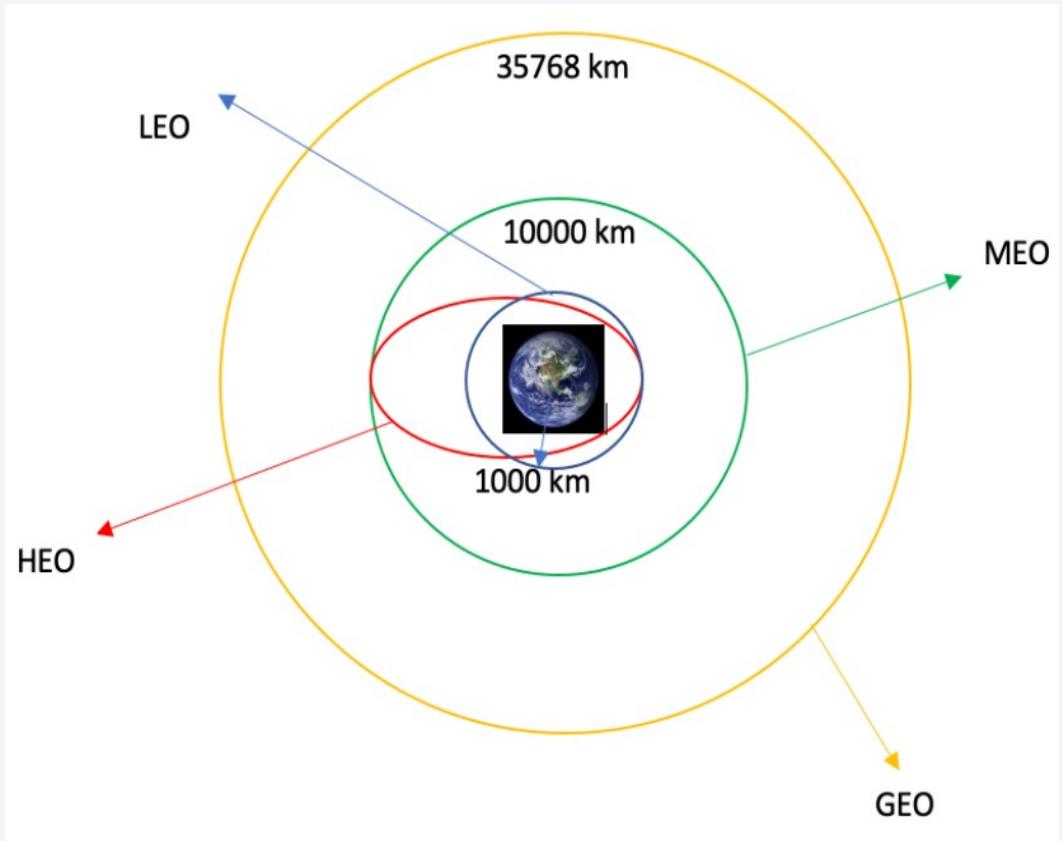
```
3. Extract all column names from the HTML table header
```

```
In [10]: column_names = []  
  
# Apply find_all() function with 'th' element on first_launch_table  
# Iterate each th element and apply the provided extract_column_from_header() to get a column name  
# Append the Non-empty column name ('if name is not None and len(name) > 0') into a list called column_names  
  
element = soup.find_all('th')  
for row in range(len(element)):  
    try:  
        name = extract_column_from_header(element[row])  
        if (name is not None and len(name) > 0):  
            column_names.append(name)  
    except:  
        pass
```

```
4. Create a dataframe by parsing the launch HTML tables
```

```
5. Export data to csv
```

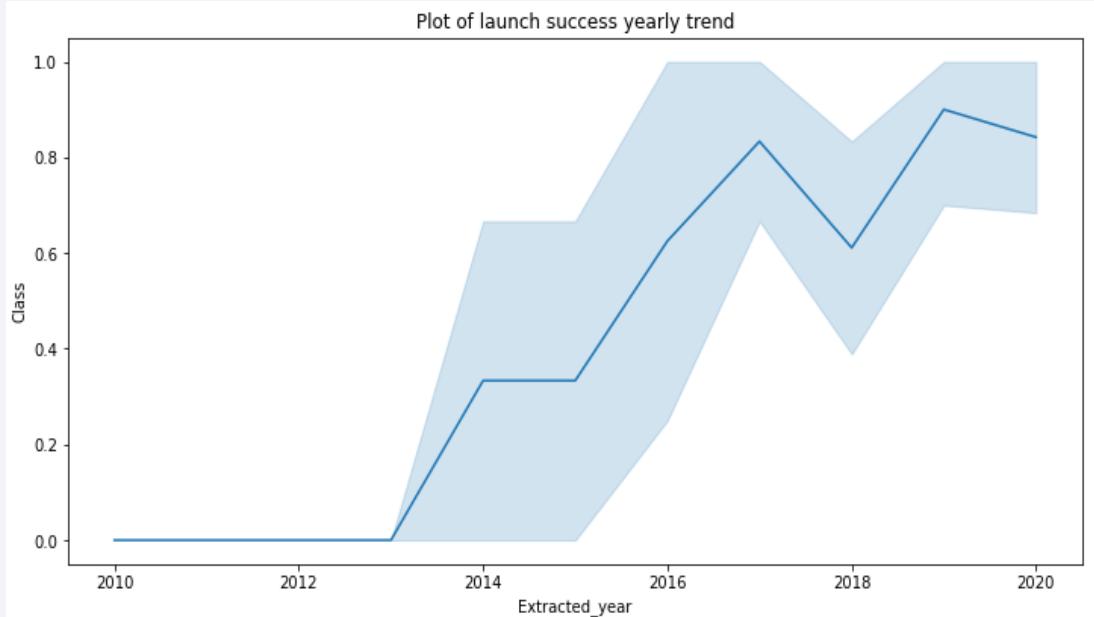
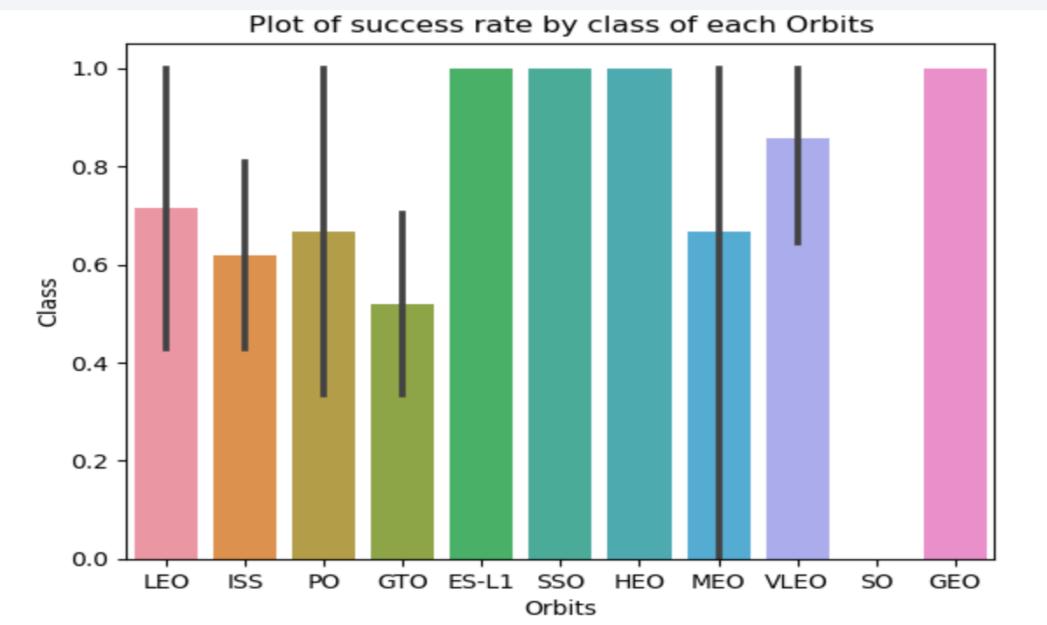
Data Wrangling



- We performed exploratory data analysis and determined the training labels.
- We calculated the number of launches at each site, and the number and occurrence of each orbits
- We created landing outcome label from outcome column and exported the results to csv.
- The link to the notebook is
[“<https://github.com/ujjwalasingh22/IBM-data-science/blob/master/Data%20wrangling.ipynb>”](https://github.com/ujjwalasingh22/IBM-data-science/blob/master/Data%20wrangling.ipynb)

EDA with Data Visualization

- We explored the data by visualizing the relationship between flight number and launch Site, payload and launch site, success rate of each orbit type, flight number and orbit type, the launch success yearly trend.



- The link to the notebook is "<https://github.com/ujjwalasingh22/IBM-data-science/blob/master/Exploratory%20Analysis%20Using%20Pandas%20and%20Matplotlib.ipynb>"

EDA with SQL

- We applied EDA with SQL to get insight from the data. We wrote queries to find out for instance:
 - 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.
- The link to the notebook is “<https://github.com/ujjwalasingh22/IBM-data-science/blob/master/Exploratory%20Analysis%20Using%20SQL.ipynb>”

Build an Interactive Map with Folium

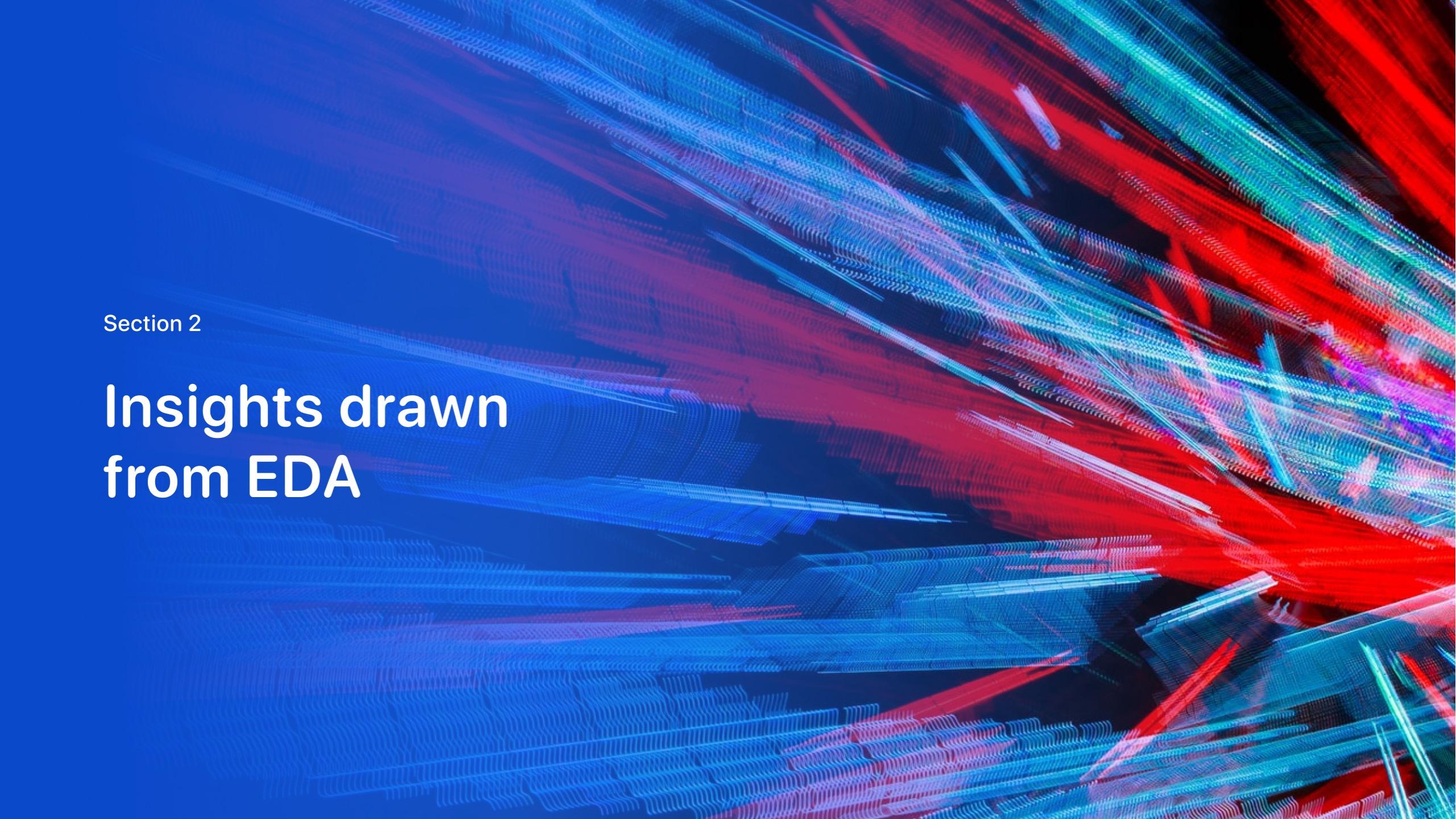
- We marked all launch sites, and added map objects such as markers, circles, lines to mark the success or failure of launches for each site on the folium map.
- We assigned the feature launch outcomes (failure or success) to class 0 and 1.i.e., 0 for failure, and 1 for success.
- Using the color-labeled 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.

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.
- The link to the notebook is “[https://github.com/ujjwalasingh22/IBM-data-science/blob/master/Predictive%20Analysis%20\(Classification\).ipynb](https://github.com/ujjwalasingh22/IBM-data-science/blob/master/Predictive%20Analysis%20(Classification).ipynb)”

Results

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

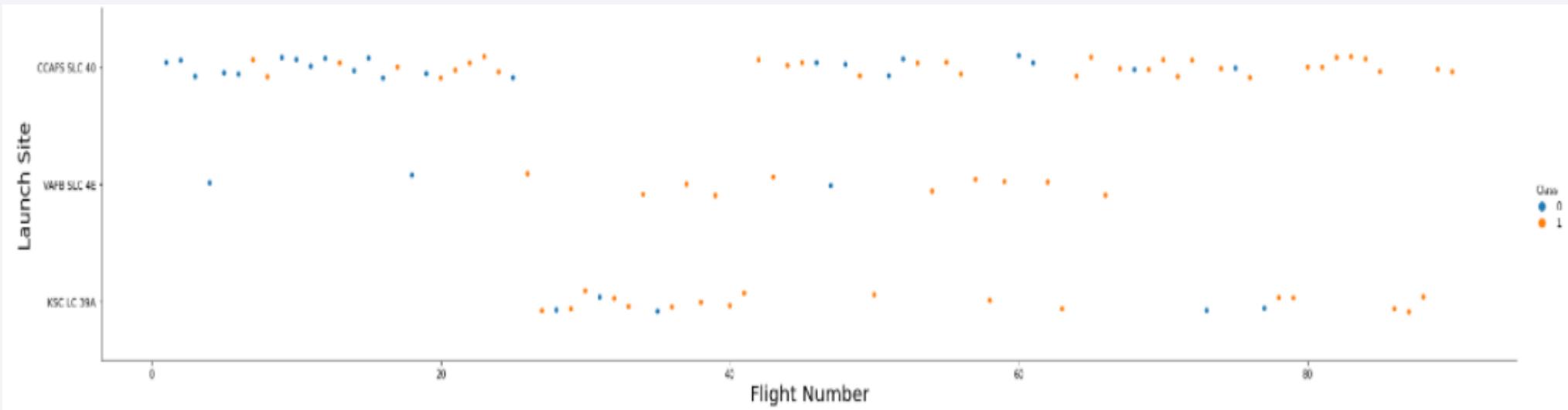
The background of the slide features a complex, abstract digital visualization. It consists of numerous thin, glowing lines that create a sense of depth and motion. The lines are primarily blue and red, with some green and purple highlights. They form a grid-like structure that curves and twists across the frame, resembling a three-dimensional space or a network of data points. The overall effect is futuristic and dynamic.

Section 2

Insights drawn from EDA

Flight Number vs. Launch Site

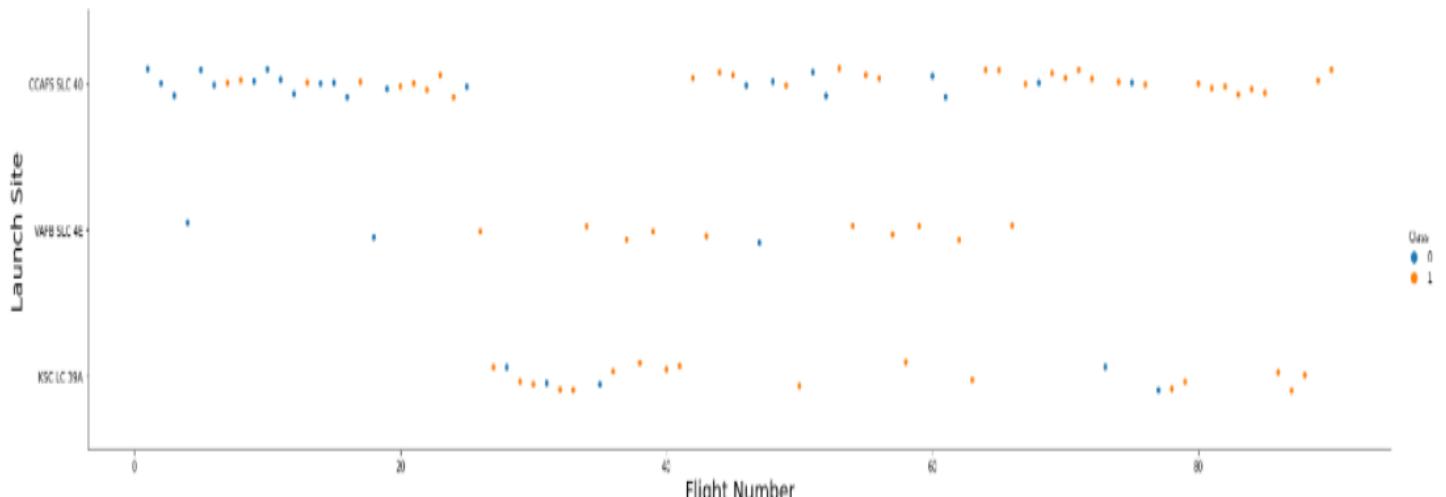
- From the plot, we found that the larger the flight amount at a launch site, the greater the success rate at a launch site.



Payload vs. Launch Site

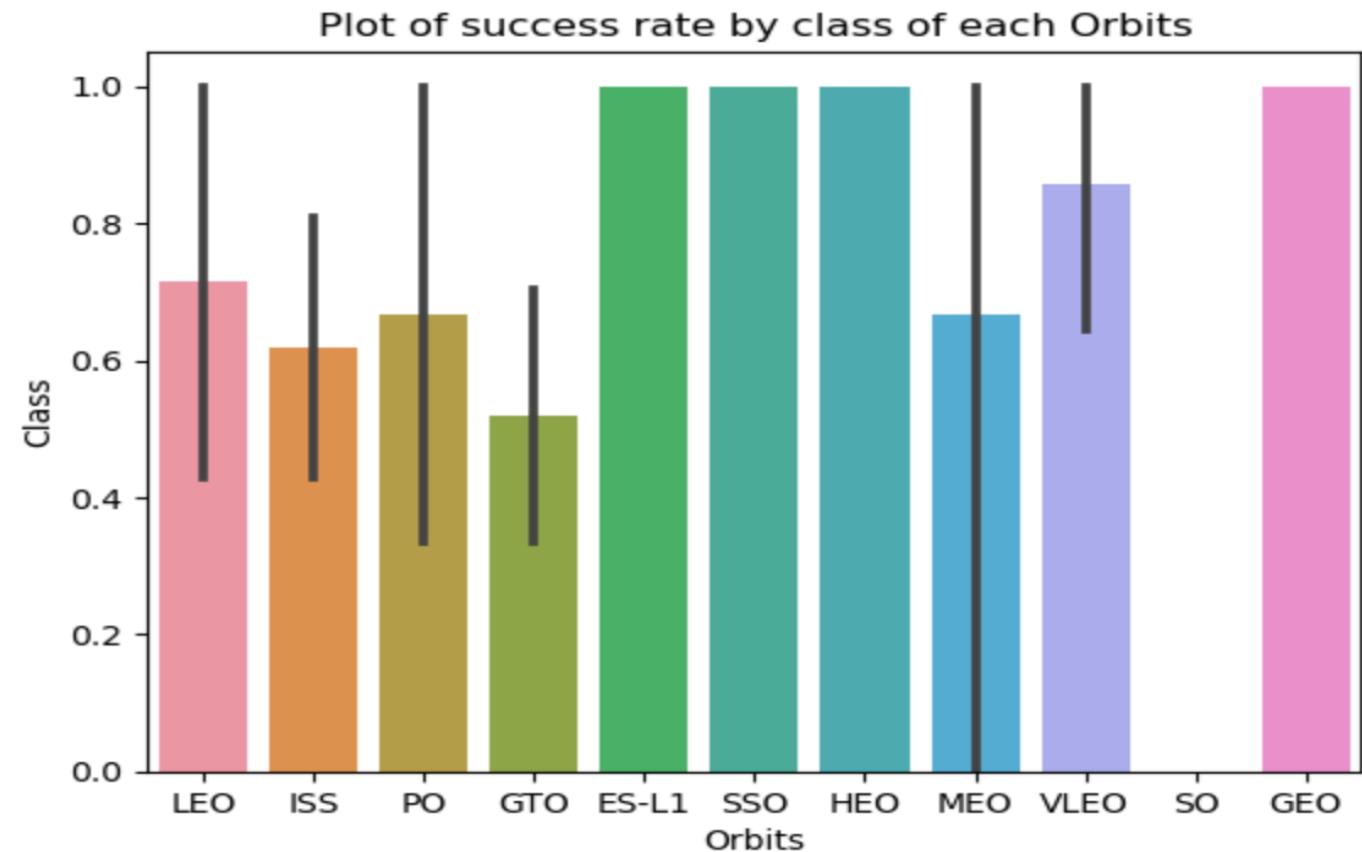


The greater the payload mass for launch site CCAFS SLC 40 the higher the success rate for the rocket.



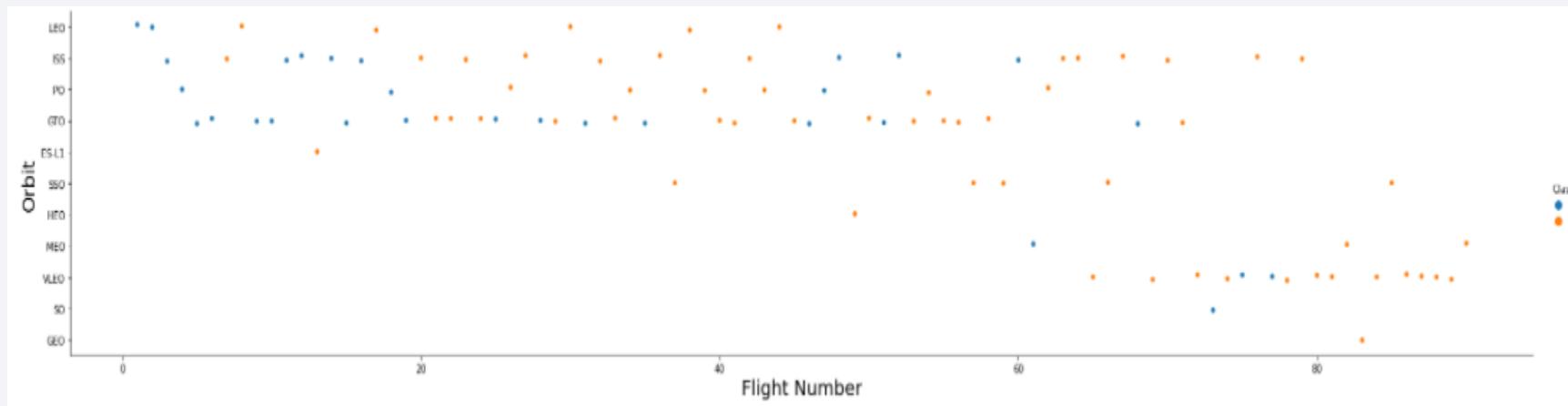
Success Rate vs. Orbit Type

- From the plot, we can see that ES-L1, GEO, HEO, SSO, VLEO had the most success rate.



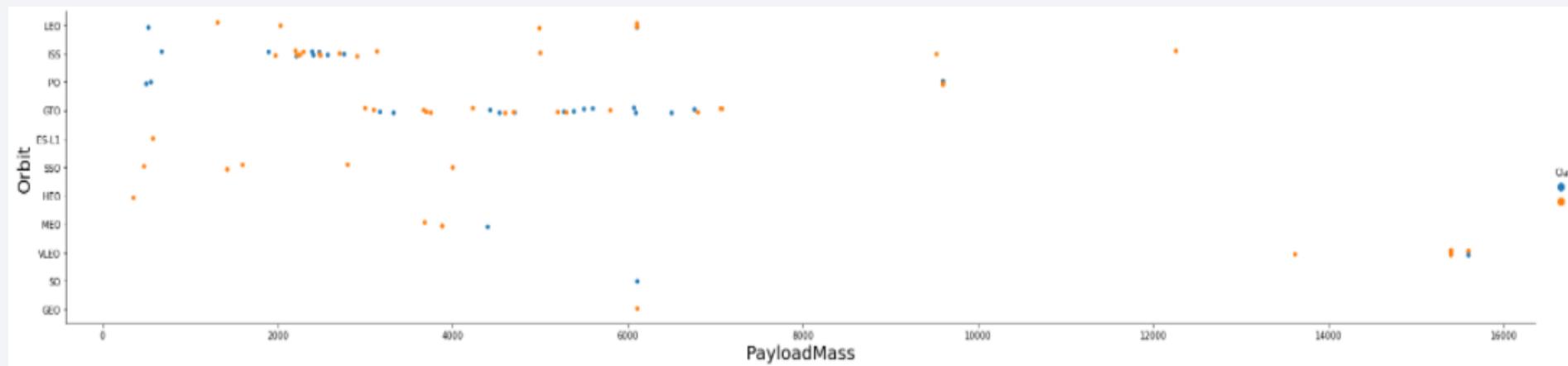
Flight Number vs. Orbit Type

- The plot below shows the Flight Number vs. Orbit type. We observe that in the 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.



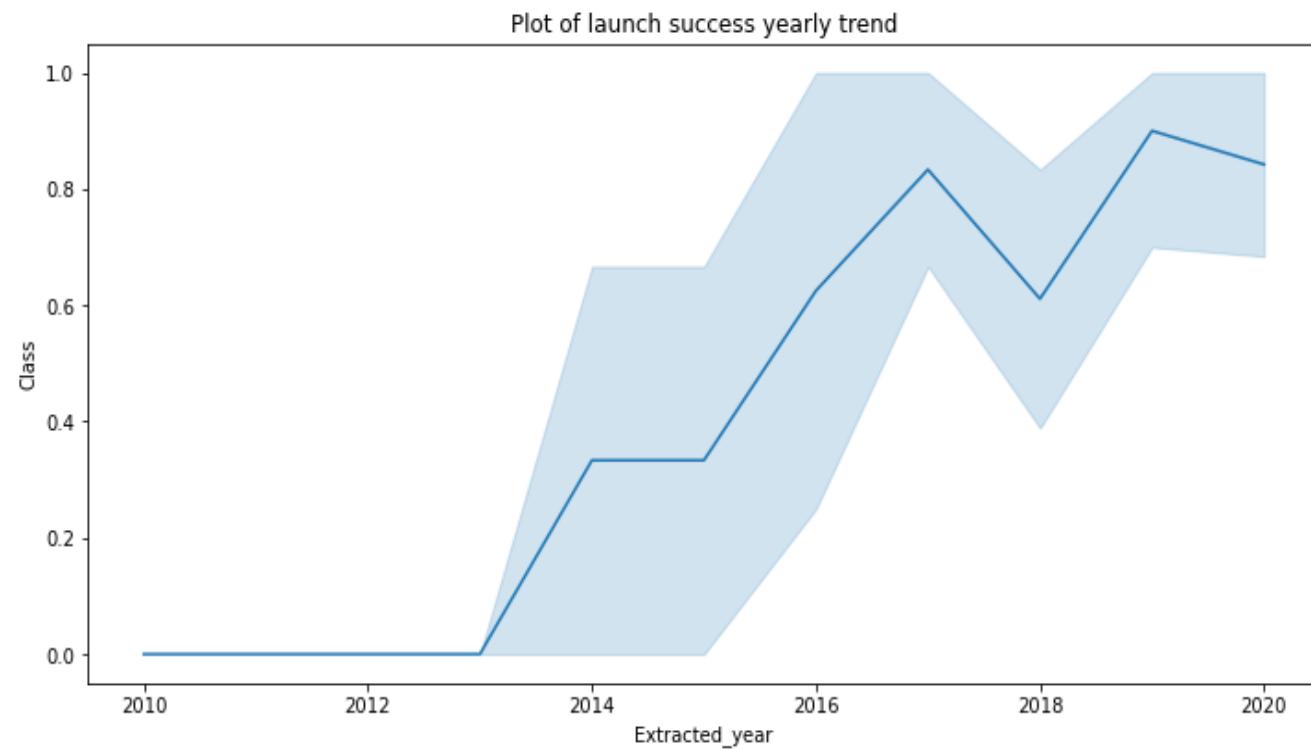
Payload vs. Orbit Type

- We can observe that with heavy payloads, the successful landing are more for PO, LEO and ISS orbits.



Launch Success Yearly Trend

- From the plot, we can observe that success rate since 2013 kept on increasing till 2020.



All Launch Site Names

- We used the key word **DISTINCT** to show only unique launch sites from the SpaceX data.

In [30]:

```
%%sql
select distinct launch_site from spacex
* ibm_db_sa://rmq23084:***@8e359033-a1c9-464
Done.
```

Out [30]: **launch_site**

CCAFS LC-40

CCAFS SLC-40

KSC LC-39A

VAFB SLC-4E

Launch Site Names Begin with 'CCA'

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

In [53]:

```
%%sql
select * from spacex where launch_site like 'CCA%' limit 5
* ibm_db_sa://rmq23084:***@e359033-a1c9-4643-82ef-8ac06f5107eb.bs2io90108kqb1od81cg.databases.appdomain.cloud:30120/bludb
Done.
```

Out[53]:

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

- We used the query above to display 5 records where launch sites begin with 'CCA'

Total Payload Mass

- We calculated the total payload carried by boosters from NASA as 45596 using the query below

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

In [60]:

```
%%sql  
select sum(payload_mass_kg_) from spacex where customer like '%NASA (CRS)'
```

```
* ibm_db_sa://rmq23084:***@8e359033-a1c9-4643-82ef-8ac06f5107eb.bs2io90108kqb1o  
Done.
```

Out[60]:

```
1  
24624
```

Average Payload Mass by F9 v1.1

- We calculated the average payload mass carried by booster version F9 v1.1 as 3676

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

```
%%sql
select avg(payload_mass_kg_) from spacex where booster_version = 'F9 v1.1'

* ibm_db_sa://rmq23084:***@8e359033-a1c9-4643-82ef-8ac06f5107eb.bs2io90108kq
Done.

1
3676
```

First Successful Ground Landing Date

- We observed that the dates of the first successful landing outcome on ground pad was 5th Jan 2017

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 like '%Success (ground pad)'
```

```
* ibm_db_sa://rmq23084:***@8e359033-a1c9-4643-82ef-8ac06f5107eb.bs2io90108kqb1od81  
Done.
```

1

2017-01-05

Successful Drone Ship Landing with Payload 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_ > 4000 and payload_mass_kg_ < 6000
* ibm_db_sa://rmq23084:***@8e359033-a1c9-4643-82ef-8ac06f5107eb.bs2io90108kqb1od81cg.databases.appdomain.cloud:30120/bludb
Done.
booster_version
F9 FT B1022
F9 FT B1031.2
```

- 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 greater than 4000 but less than 6000

Total Number of Successful and Failure Mission Outcomes

List the total number of successful and failure mission outcomes

```
%%sql  
select count(mission_outcome) as success_outcome from spacex where mission_outcome like '%Success%'
```

```
* ibm_db_sa://rmq23084:***@8e359033-alc9-4643-82ef-8ac06f5107eb.bs2io90108kqb1od8lcg.databases.appdomain.  
Done.
```

success_outcome

```
45
```

```
%%sql  
select count(mission_outcome) as failure_outcome from spacex where mission_outcome like '%Failure%'
```

```
* ibm_db_sa://rmq23084:***@8e359033-alc9-4643-82ef-8ac06f5107eb.bs2io90108kqb1od8lcg.databases.appdomain.  
Done.
```

failure_outcome

```
0
```

- We used wildcard like '%' to filter for **WHERE** MissionOutcome was a success or a failure.

Boosters Carried Maximum Payload

- We determined the booster that have carried the maximum payload using a subquery in the WHERE clause and the MAX() function.

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

```
%%sql
select booster_version, payload_mass_kg_ from spacex where payload_mass_kg_ = (
    select max(payload_mass_kg_) from spacex) order by booster_version

* ibm_db_sa://rmq23084:***@8e359033-a1c9-4643-82ef-8ac06f5107eb.bs2io90108kqb1od81cg.dat
Done.

booster_version payload_mass_kg_
F9 B5 B1048.4      15600
F9 B5 B1049.4      15600
F9 B5 B1049.5      15600
F9 B5 B1058.3      15600
F9 B5 B1060.2      15600
```

2015 Launch Records

- We used combinations of the WHERE clause, LIKE, AND, and BETWEEN conditions to filter for failed landing outcomes in drone ship, their booster versions, and launch site names for year 2015

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

```
%%sql
select booster_version, launch_site, landing_outcome from spacex where landing_outcome
like '%Failure (drone ship)' and date between '2015-01-01' AND '2015-12-31'
```

```
* ibm_db_sa://rmq23084:***@8e359033-a1c9-4643-82ef-8ac06f5107eb.bs2io90108kqb1od81cg.databases.
Done.
```

booster_version	launch_site	landing_outcome
F9 v1.1 B1012	CCAFS LC-40	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

```
%%sql
select landing_outcome, count(landing_outcome) from spacex where date between '2010-06-04' AND '2017-03-20'

SELECT LandingOutcome, COUNT(LandingOutcome)
  FROM SpaceX
 WHERE DATE BETWEEN '2010-06-04' AND '2017-03-20'
 GROUP BY LandingOutcome
 ORDER BY COUNT(LandingOutcome) DESC
```

- We selected Landing outcomes and the COUNT of landing outcomes from the data and used the WHERE clause to filter for landing outcomes BETWEEN 2010-06-04 to 2010-03-20.
- We applied the GROUP BY clause to group the landing outcomes and the ORDER BY clause to order the grouped landing outcome in descending order.

The background of the slide is a photograph taken from space at night. It shows the curvature of the Earth against a dark blue-black void of space. City lights are visible as numerous small white and yellow dots, primarily concentrated in the lower right quadrant where the United States appears. In the upper right, there are bright green and yellow bands of light, likely the Aurora Borealis or Australis. The overall atmosphere is dark and mysterious.

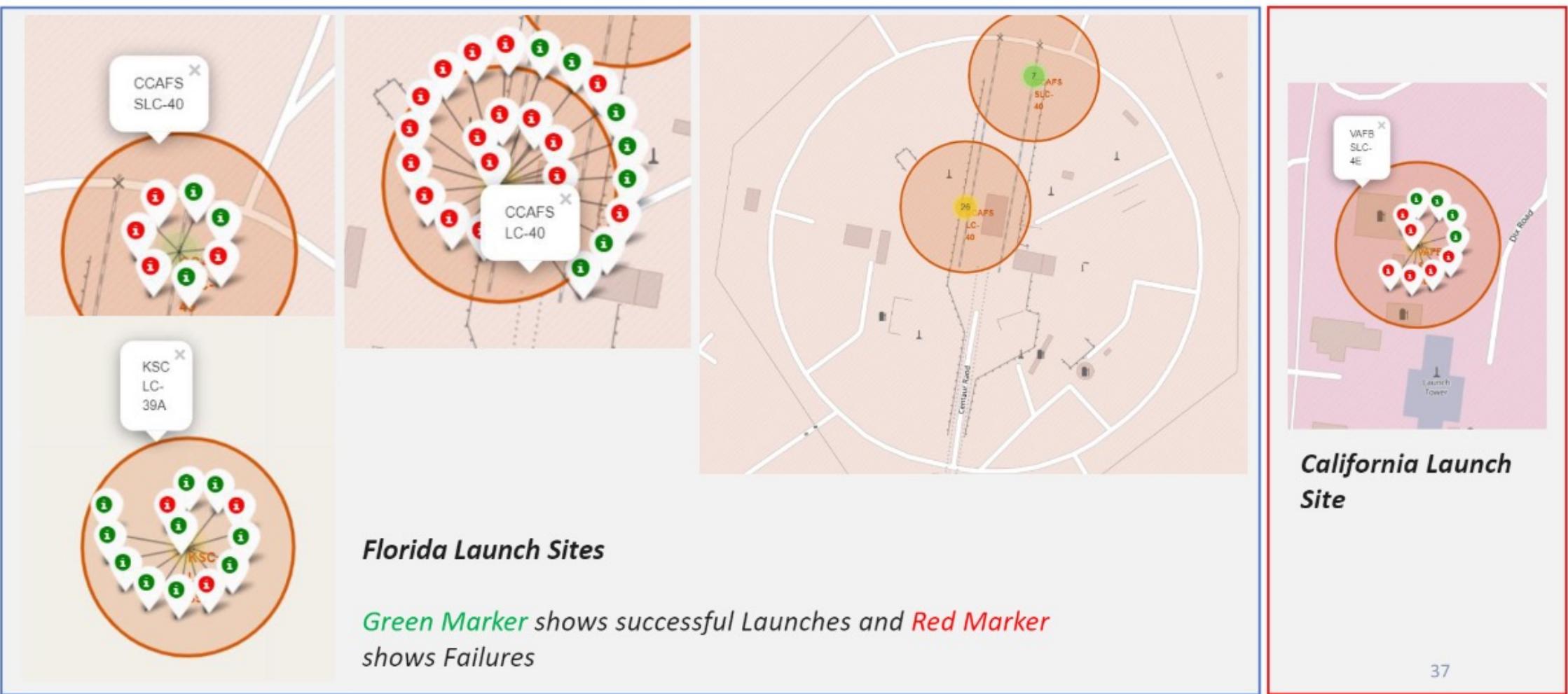
Section 4

Launch Sites Proximities Analysis

All launch sites global map markers



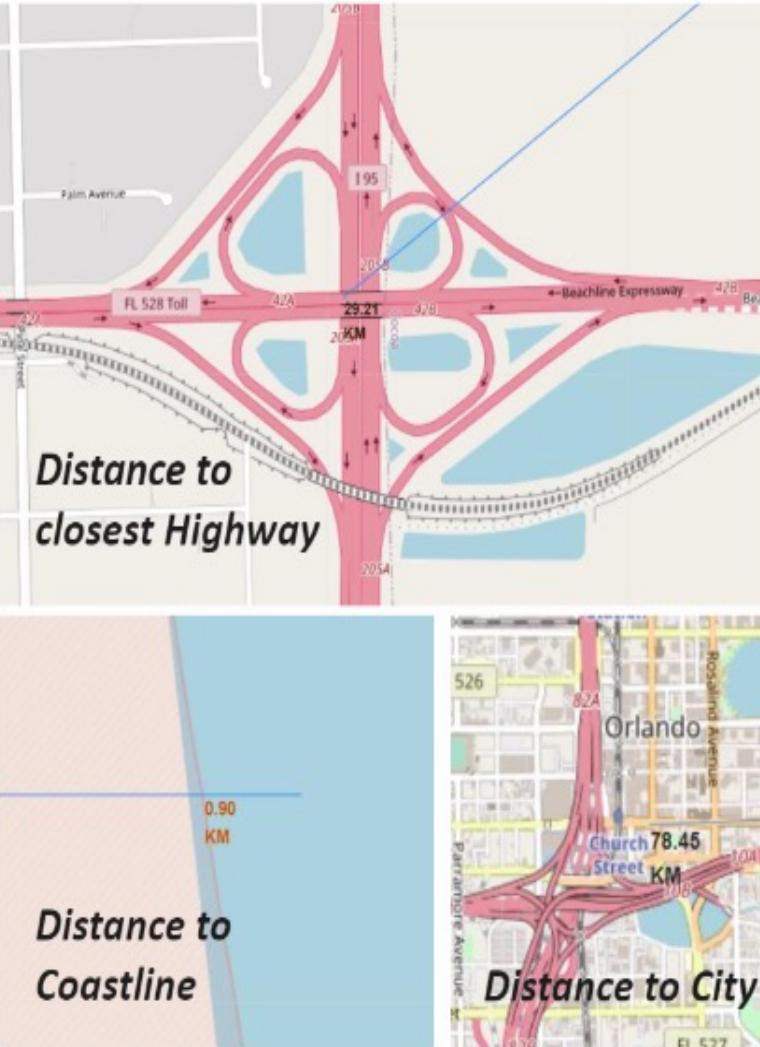
Markers showing launch sites with color labels



37

35

Launch Site distance to landmarks



- 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 6

Predictive Analysis (Classification)

Classification Accuracy

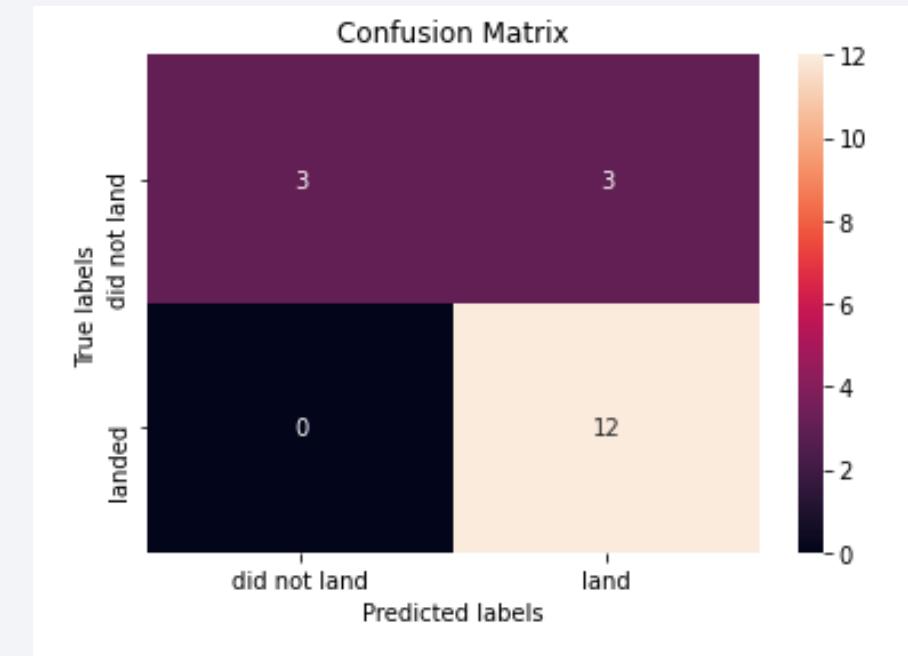
- The decision tree classifier is the model with the highest classification accuracy

```
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

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!

