

Winning Space Race with Data Science

Andrew Yeboah Nortey
April 2024



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 that need answers

1. Which factors are behind the failure of landing?
2. Will the rockets land successfully?
3. What the accuracy of a successful landing?

Section 1

Methodology

Methodology

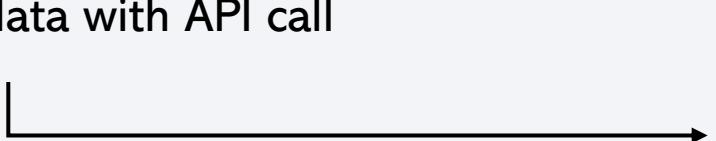
Executive Summary

- Data collection methodology:
 - With Rest API and Web Scrapping
- Perform data wrangling
 - Data were transformed and one hot encoded to be apply later on the Machine Learning models
- Perform exploratory data analysis (EDA) using visualization and SQL
 - Discovering new patterns in the data with visualization techniques such as scatter plots
- Perform interactive visual analytics using Folium and Plotly Dash
 - Dash and Folium were used to achieve this goal
- Perform predictive analysis using classification models
 - Classification machine learning models were built to achieve this goal

Data Collection

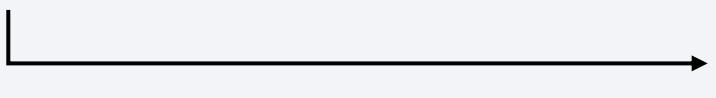
Data sets were collected using the API call from several websites, I collected rocket, launchpad, payloads, and cores data from <https://api.spacexdata.com/v4> website.

1. Collecting the data with API call



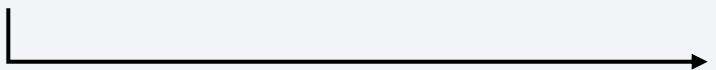
2. Converting to data frame with help of JSON

3. Updating columns and rows (pre-processing)



4. Filtering the data to keep only Falcon 9 launches

5. Convert the data to csv file with name
'dataset_part_1.csv'



Data Collection – SpaceX API

1. Collecting the data with API call

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

3. Updating columns and rows (pre-processing)

```
# Lets take a subset of our dataframe keeping only the features we want and the flight_number..and date_utc.
data = data[['rocket', 'payloads', 'launchpad', 'cores', 'flight_number', 'date_utc']]
# We will remove rows with multiple cores because those are falcon rockets with 2 extra rocket boosters and rows that have multiple payloads in a single
# core
data = data[data['cores'].map(len)==1]
data = data[data['payloads'].map(len)==1]
# Since payloads and cores are lists of size 1 we will also extract the single value in the list and replace the feature.
data['cores'] = data['cores'].map(lambda x: x[0])
data['payloads'] = data['payloads'].map(lambda x: x[0])
# We also want to convert the date_utc to a datetime datatype and then extracting the date leaving the time
data['date'] = pd.to_datetime(data['date_utc']).dt.date
# Using the date we will restrict the dates of the launches
data = data[data['date'] <= datetime.date(2020, 11, 13)]
```

5. Convert the data to csv file with name 'dataset_part_1.csv'

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

2. Converting to data frame with help of JSON

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

Using the dataframe data print the first 5 rows

```
# Get the head of the dataframe
data.head(5)
```

	static_fire_date_utc	static_fire_date_unix	tbd	net	window	rocket	success	details	crew	ships	capsules	payloads
0	2006-03-17T00:00:00.000Z	1.142554e+09	False	False	0.0	5e9d0d95eda69955f709d1eb	False	Engine failure at 33 seconds and loss of vehicle	0	0	0	[5eb0e4b5b6c3bb0006eeb1e1]

4. Filtering the data to keep only Falcon 9 launches

```
# Hint data['BoosterVersion']!='Falcon 1'
data_falcon9 = data2[data2['BoosterVersion']!='Falcon 1']
```

GitHub repo: <https://github.com/aynortey/IBM-Data-Science-Capstone-Project/blob/main/jupyter-labs-spacex-data-collection-api1.ipynb>

Data Collection - Scraping

GitHub repo: <https://github.com/aynortey/IBM-Data-Science-Capstone-Project/blob/main/jupyter-labs-webscraping1.ipynb>

1. Creating the BeautifulSoup object

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

3. Creating the launch_dict

```
# Customer
# TODO: Append the customer into launch_dict with key `Customer`
if row[6].a != None:
    customer = row[6].a.string
else:
    customer = ''
#print(customer)
launch_dict['Customer'].append(customer)

# Launch outcome
# TODO: Append the launch_outcome into launch_dict with key `Launch_outcome`
launch_outcome = list(row[7].strings)[0]
#print(launch_outcome)
launch_dict['Launch outcome'].append(launch_outcome)

# Booster landing
# TODO: Append the booster_landing into launch_dict with key `Booster_landing`
booster_landing = landing_status(row[8])
#print(booster_landing)
launch_dict['Booster landing'].append(booster_landing)
```

5. Convert the data to csv file with name 'spacex_web_scraped.csv'

```
df.to_csv('spacex_web_scraped.csv', index=False)
```

2. Getting column names

```
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
tc = first_launch_table.find_all('th')
for th in tc:
    name = extract_column_from_header(th)
    if name is not None and len(name) > 0:
        column_names.append(name)
```

4. Converting to final data frame

Launch Dict										
	Flight No.	Launch site	Payload	Payload mass	Orbit	Customer	Launch outcome	Version Booster	Booster landing	Date Time
0	1	CCAFS	Dragon Spacecraft Qualification Unit	0	LEO	SpaceX	Success\n	F9 v1.0B0003.1	Failure	4 June 2010 18:45
1	2	CCAFS	Dragon	0	LEO	NASA	Success	F9 v1.0B0004.1	Failure	8 December 2010 15:43
2	3	CCAFS	Dragon	525 kg	LEO	NASA	Success	F9 v1.0B0005.1	No attempt\n	22 May 2012 07:44
3	4	CCAFS	SpaceX CRS-1	4,700 kg	LEO	NASA	Success\n	F9 v1.0B0006.1	No attempt	8 October 2012 00:35
4	5	CCAFS	SpaceX CRS-2	4,877 kg	LEO	NASA	Success\n	F9 v1.0B0007.1	No attempt\n	1 March 2013 15:10

Data Wrangling

GitHub repo: <https://github.com/aynortey/IBM-Data-Science-Capstone-Project/blob/main/labs-jupyter-spacex-Data%20wrangling1.ipynb>

1. Loading the data set

```
df=pd.read_csv("https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBM-DS0321EN-SkillsNetwork/datasets/dataset_part_1.csv")
df.head(10)
```

3. Finding the bad outcomes

We create a set of outcomes where the second stage did not land successfully:

```
bad_outcomes=set(landing_outcomes.keys()[[1,3,5,6,7]])
bad_outcomes

{'False ASDS', 'False Ocean', 'False RTLS', 'None ASDS', 'None None'}
```

5. Determining the success outcome

```
df["Class"].mean()
```

```
0.6666666666666666
```

6. Convert the data to csv file with name 'dataset_part_2.csv'

```
df.to_csv("dataset_part_2.csv", index=False)
```

2. Creating landing outcomes

```
# landing_outcomes = values on Outcome column
landing_outcomes = df["Outcome"].value_counts()
landing_outcomes
```

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

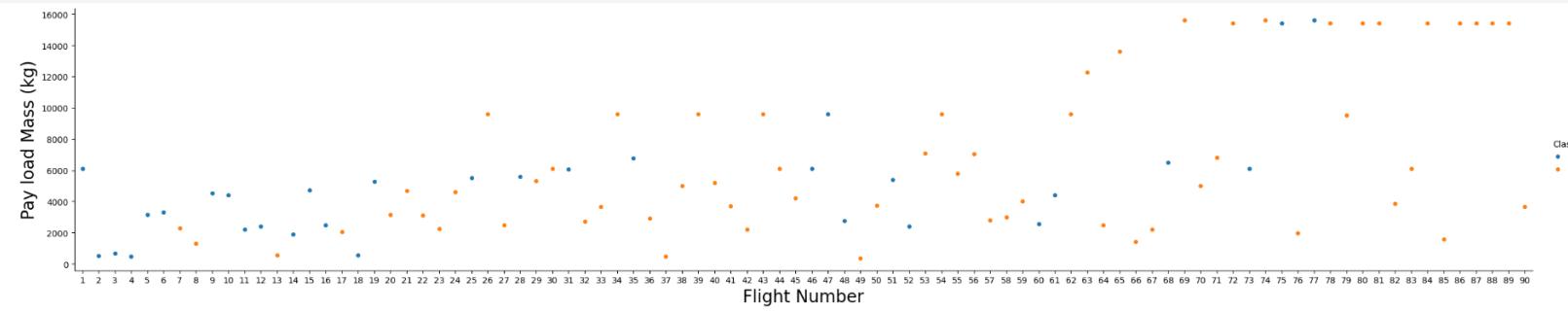
4. Presenting outcomes as 0 and 1

```
df['Class']=landing_class
df[['Class']].head(8)
```

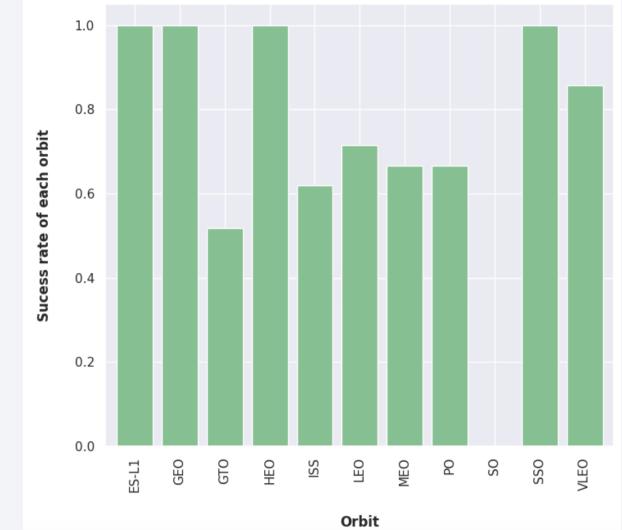
Class	
0	0
1	0
2	0
3	0
4	0
5	0
6	1
7	1

EDA with Data Visualization

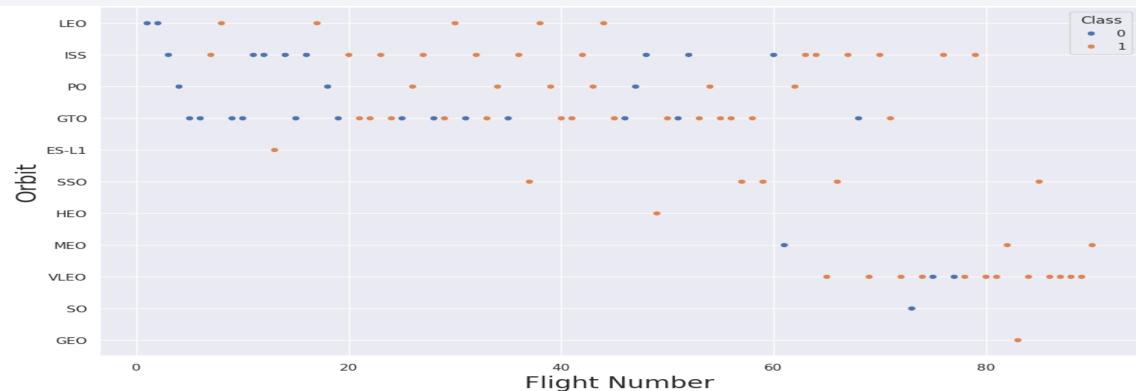
GitHub repo: <https://github.com/aynortey/IBM-Data-Science-Capstone-Project/blob/main/jupyter-labs-eda-dataviz-jupyterlite1.ipynb>



Categorial plot between Flight number and Pay load mass (kg)

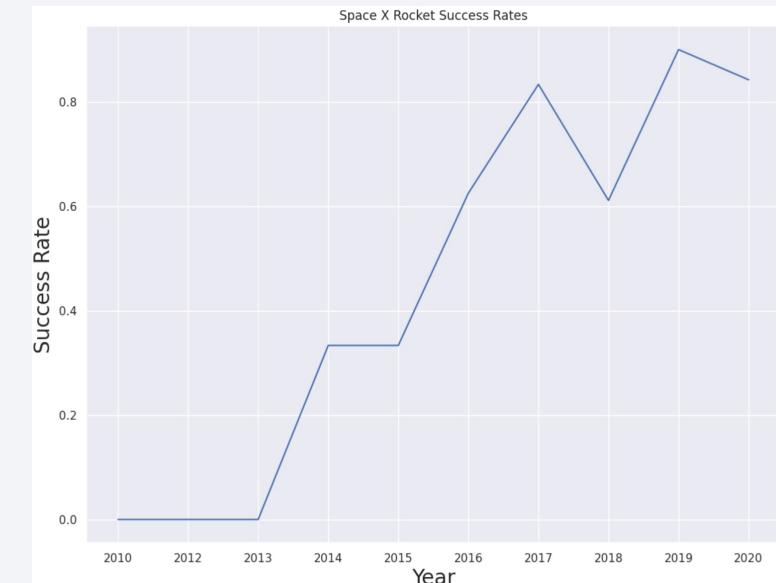


Bar chart between Orbit and Success rate of each orbit



Line chart between Year and Success rate

Scatter plot between Orbit and Flight number



I used SQL queries to answer the following questions:

- 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 (CRS)
- 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. Use a subquery
- List the failed landing_outcomes in drone ship, their booster versions, and launch site names for the 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

Build an Interactive Map with Folium

- folium.Marker() was used to create marks on the maps.
- folium.Circle() was used to create a circles above markers on the map.
- folium.Icon() was used to create an icon on the map.
- folium.PolyLine() was used to create polynomial line between the points.
- folium.plugins.AntPath() was used to create animated line between the points.
- markerCluster() was used to simplify the maps which contain several markers with identical coordination.

GitHub repo: <https://github.com/aynorsey/IBM-Data-Science-Capstone-Project/blob/main/Interactive%20Visual%20Analytics%20with%20Folium%20lab.ipynb>

Build a Dashboard with Plotly Dash

- Dash and html components were used as they are the most important thing and almost everything depends on them, such as graphs, tables, dropdowns, etc.
- Pandas was used to simplifying the work by creating dataframe.
- Plotly was used to plot the graphs.
- Pie chart and scatter chart were used to for plotting purposes.
- RangeSlider was used for payload mass range selection.
- Dropdown was used for launch sites.

Predictive Analysis (Classification)

1. Building the model

Create column for the class

Standardize the data

Split the data info train and test sets

Build GridSearchCV model and fit the data

3. Finding the optimal model

Find the best hyperparameters for the models

Find the best model with highest accuracy

Confirm the optimal model

2. Evaluating the model

Calculating the accuracies

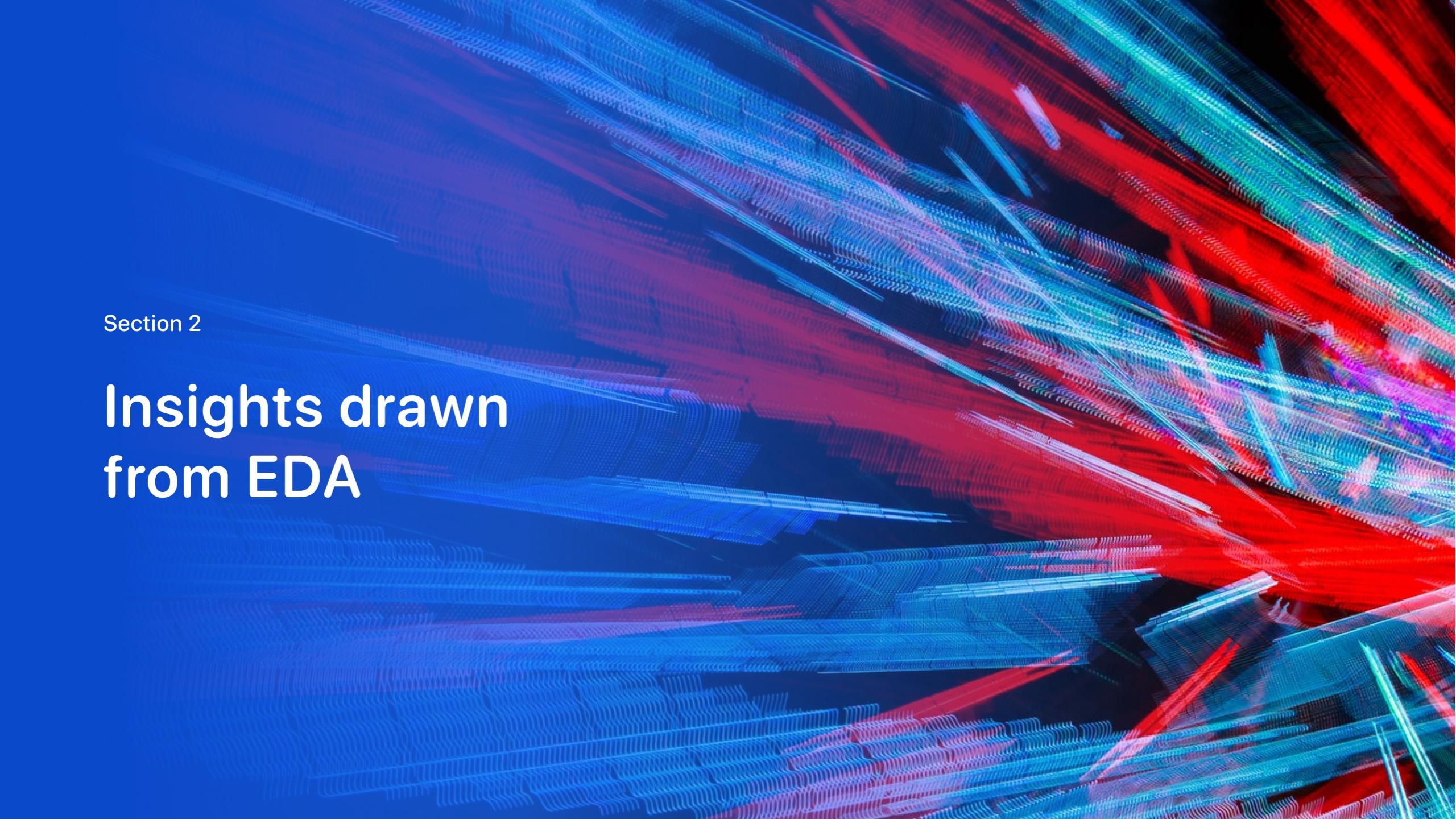
Calculating the confusion matrixes

Plot the results

GitHub repo: https://github.com/aynortey/IBM-Data-Science-Capstone-Project/blob/main/SpaceX_Machine%20Learning%20Prediction_Part_5.ipynb

Results

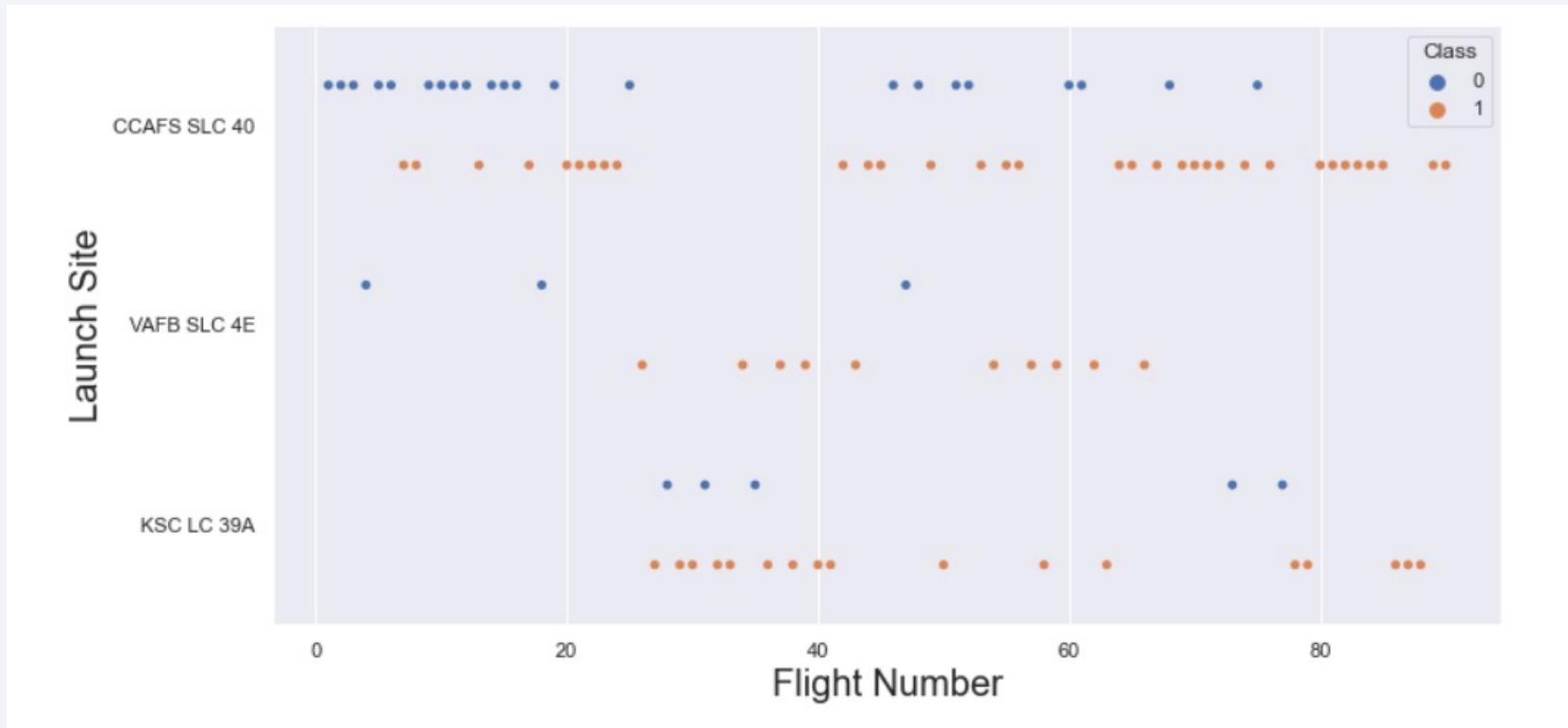
- Exploratory data analysis results
- Interactive analytics demo in screenshots
- Predictive analysis results

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

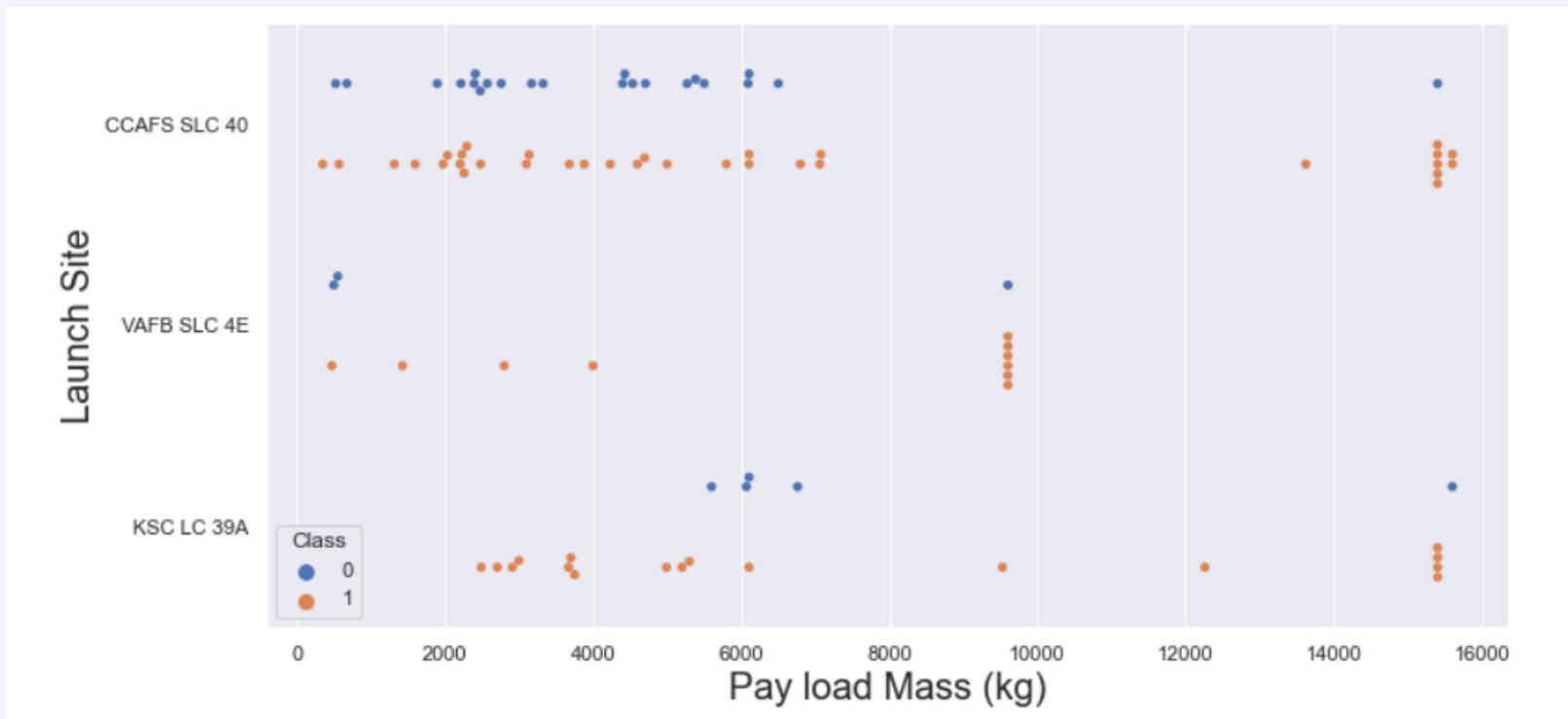
Flight Number vs. Launch Site



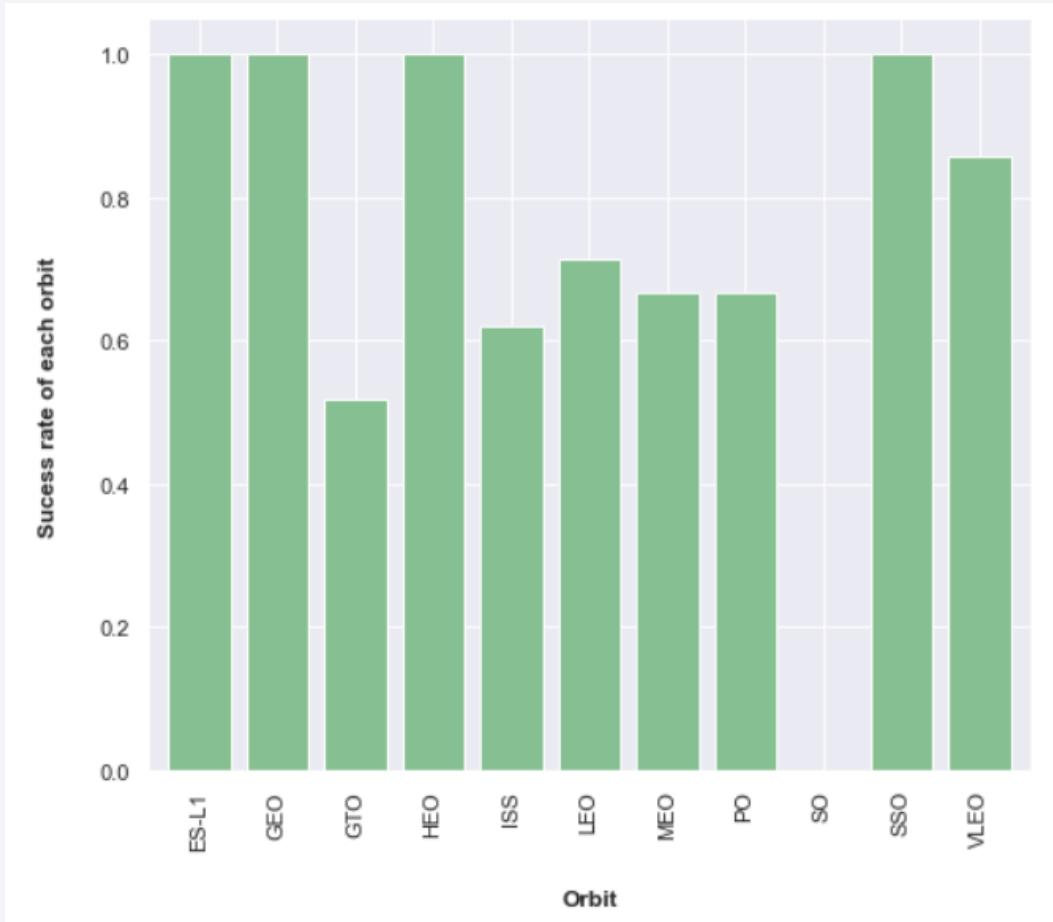
With the increase of flight number, the success rate is increasing as well in the launch sites

Payload vs. Launch Site

With the increase of Pay load Mass, the success rate is increasing as well in the launch sites



Success Rate vs. Orbit Type

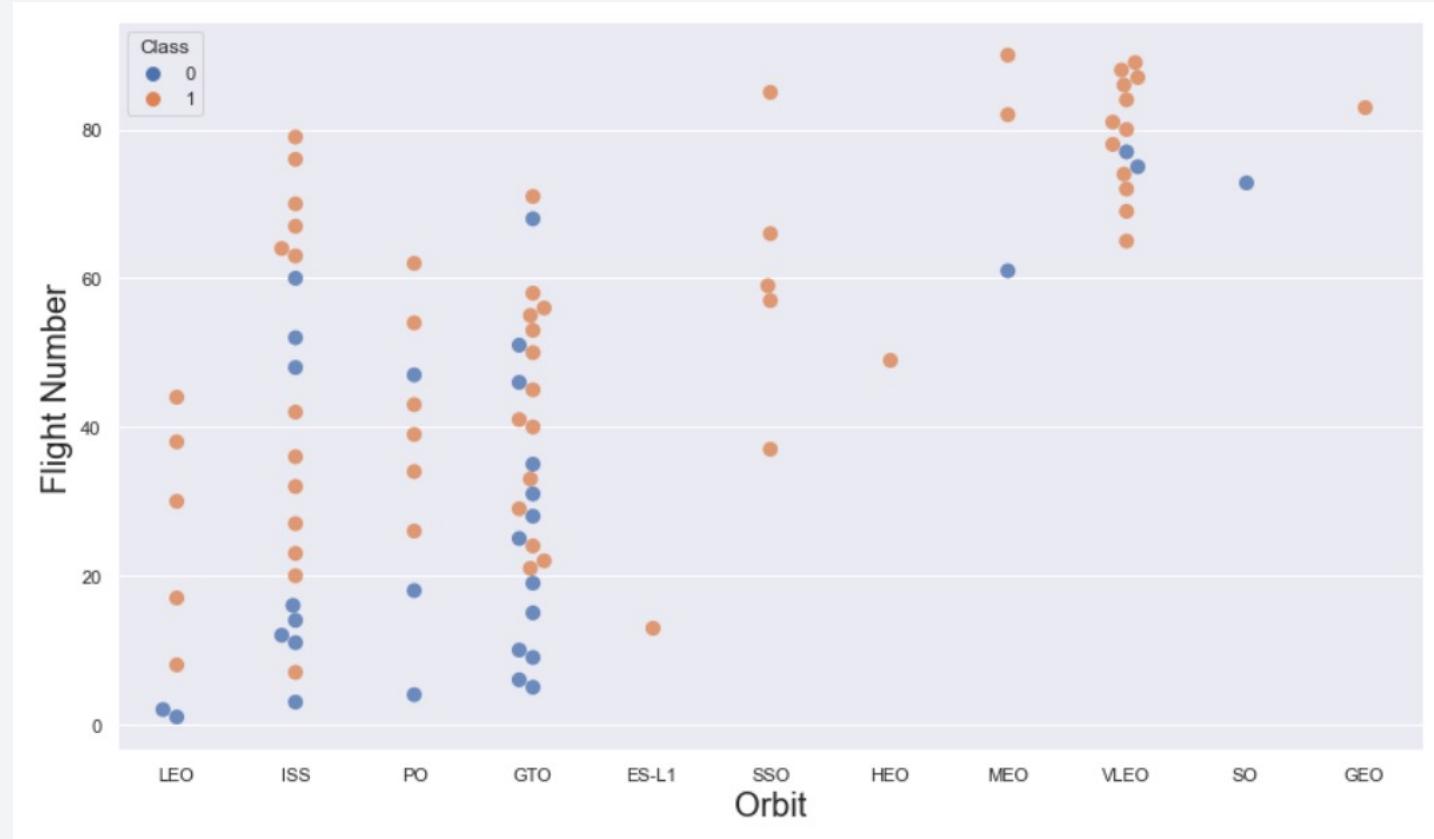


ES-L1, GEO, HEO, and SSO have a success rate of 100%

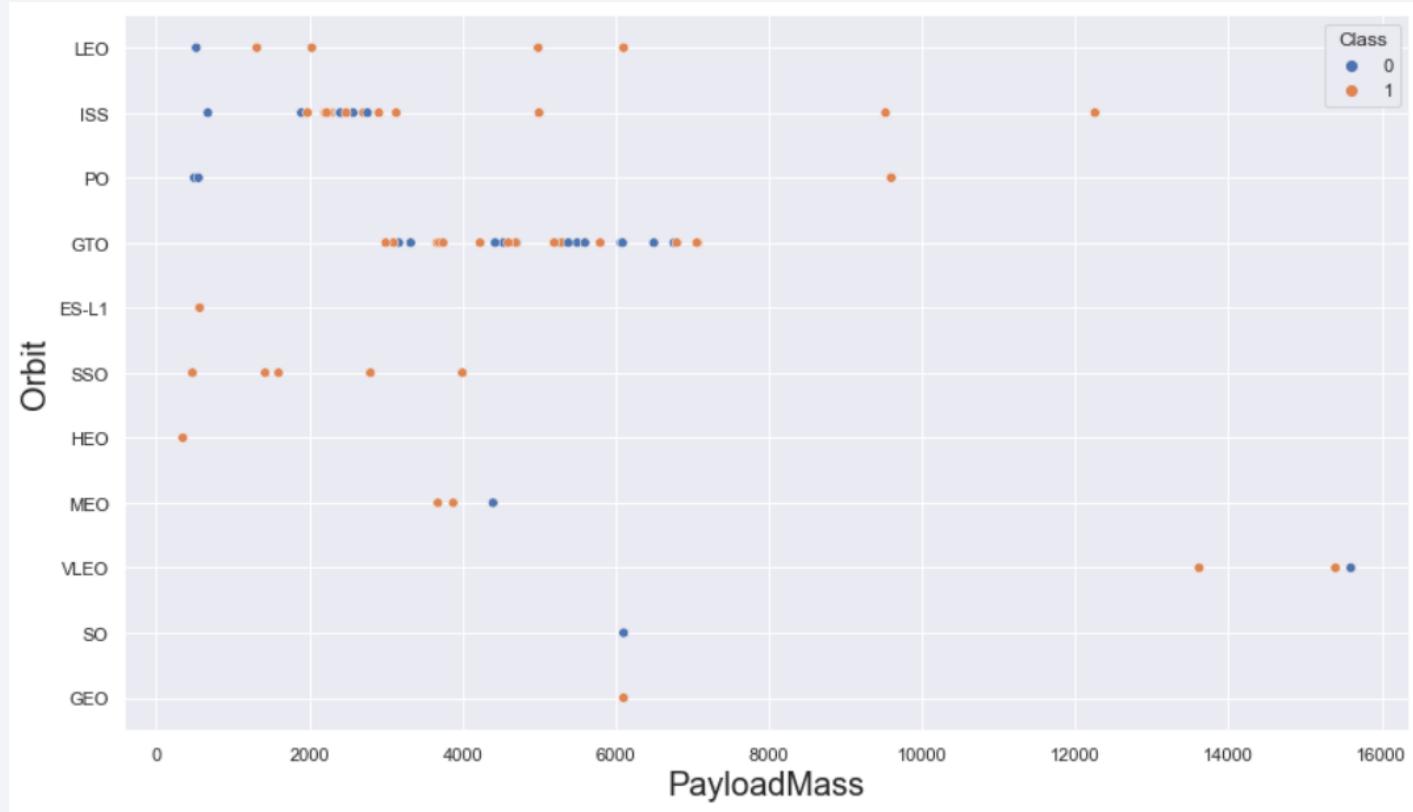
SO has a success rate of 0%

Flight Number vs. Orbit Type

It's hard to tell anything here, but we can say there is no actual relationship between flight number and GTO.



Payload vs. Orbit Type

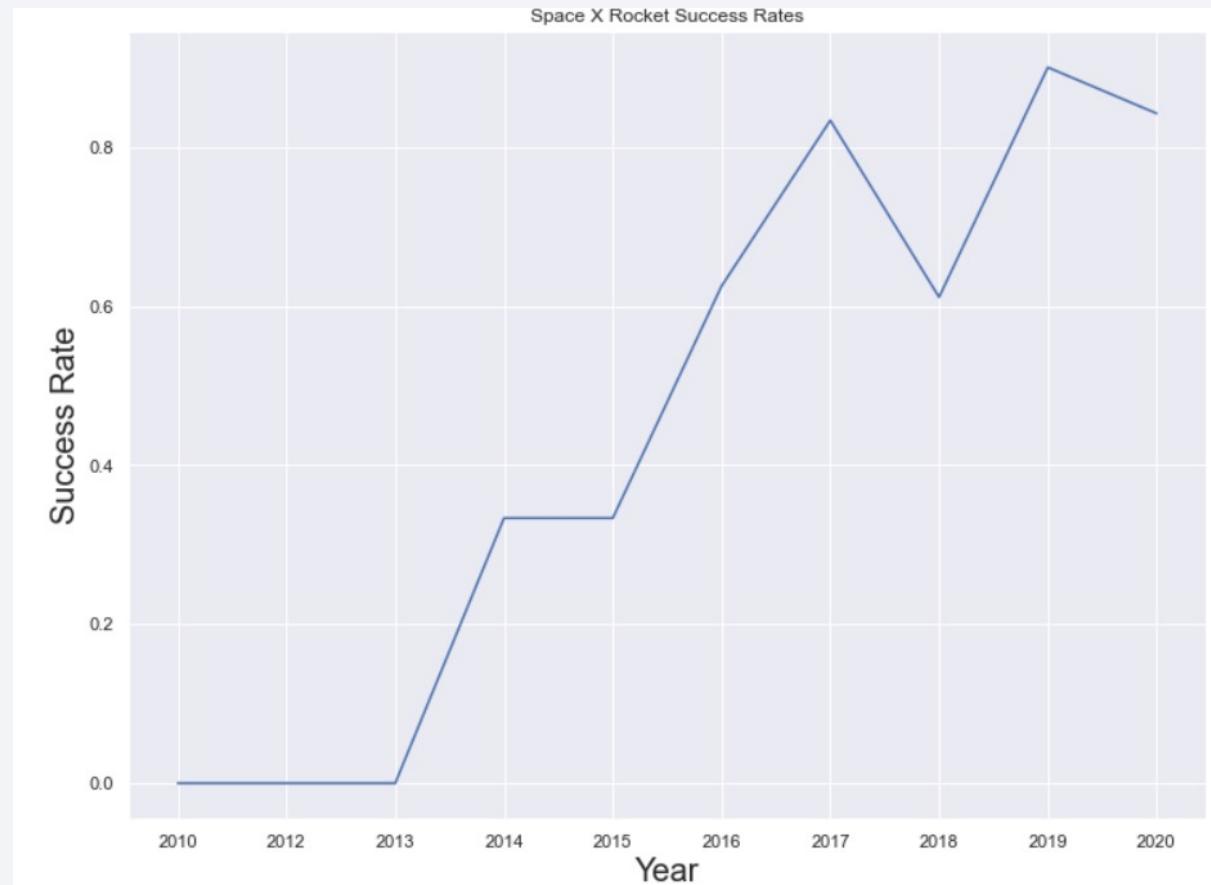


First thing to see is how the Pay load Mass between 2000 and 3000 is affecting ISS.

Similarly, Pay load Mass between 3000 and 7000 is affecting GTO.

Launch Success Yearly Trend

Since the year 2013, there was a massive increase in success rate. However, it dropped little in 2018 but later it got stronger than before.



All Launch Site Names

```
%sql SELECT DISTINCT LAUNCH_SITE as "Launch_Sites" FROM SPACEX;
```

Launch_Sites
CCAFS LC-40
CCAFS SLC-40
KSC LC-39A
VAFB SLC-4E

We can get the unique values by using “DISTINCT”

Launch Site Names Begin with 'CCA'

```
%sql SELECT * FROM SPACEX WHERE LAUNCH_SITE LIKE 'CCA%' LIMIT 5;
```

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

We can get only 5 rows by using “LIMIT”

Total Payload Mass

```
%sql SELECT SUM(PAYLOAD_MASS__KG_) AS "Total Payload Mass by NASA (CRS)" FROM SPACEX WHERE CUSTOMER = 'NASA (CRS)';
```

Total Payload Mass by NASA (CRS)
45596

We can get the sum of all values by using “SUM”

Average Payload Mass by F9 v1.1

```
%sql SELECT AVG(PAYLOAD_MASS__KG_) AS "Average Payload Mass by Booster Version F9 v1.1" FROM SPACEX \
WHERE BOOSTER_VERSION = 'F9 v1.1';
```

Average Payload Mass by Booster Version F9 v1.1

2928

We can get the average of all values by using “AVG”

First Successful Ground Landing Date

```
%sql SELECT MIN(DATE) AS "First Sucessful Landing Outcome in Ground Pad" FROM SPACEX \
WHERE LANDING__OUTCOME = 'Success (ground pad)';
```

First Sucessful Landing Outcome in Ground Pad
2015-12-22

We can get the first successful data by using “MIN”, because first date is same with the minimum date

Successful Drone Ship Landing with Payload between 4000 and 6000

```
%sql SELECT BOOSTER_VERSION FROM SPACEX WHERE LANDING_OUTCOME = 'Success (drone ship)' \
AND PAYLOAD_MASS_KG_ > 4000 AND PAYLOAD_MASS_KG_ < 6000;
```

booster_version
F9 FT B1022
F9 FT B1026
F9 FT B1021.2
F9 FT B1031.2

The payload mass data was taken between 4000 and 6000 only, and the landing outcome was determined to be “success drone ship”

Total Number of Successful and Failure Mission Outcomes

```
%sql SELECT COUNT(MISSION_OUTCOME) AS "Successful Mission" FROM SPACEX WHERE MISSION_OUTCOME LIKE 'Success%';
```

Successful Mission
100

We can get the number of all the successful mission by using “COUNT” and LIKE “Success%”

```
%sql SELECT COUNT(MISSION_OUTCOME) AS "Failure Mission" FROM SPACEX WHERE MISSION_OUTCOME LIKE 'Failure%';
```

Failure Mission
1

We can get the number of all the failure mission by using “COUNT” and LIKE “Failure%”

Boosters Carried Maximum Payload

```
%sql SELECT DISTINCT BOOSTER_VERSION AS "Booster Versions which carried the Maximum Payload Mass" FROM SPACEX \
WHERE PAYLOAD_MASS__KG_ =(SELECT MAX(PAYLOAD_MASS__KG_) FROM SPACEX);
```

Booster Versions which carried the Maximum Payload Mass
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

We can get the maximum payload masses by using “MAX”

2015 Launch Records

```
%sql SELECT month(DATE) as Month, BOOSTER_VERSION, LAUNCH_SITE FROM SPACEX WHERE year(DATE) = '2015' AND \
LANDING__OUTCOME = 'Failure (drone ship)';
```

MONTH	booster_version	launch_site
1	F9 v1.1 B1012	CCAFS LC-40
4	F9 v1.1 B1015	CCAFS LC-40

We can get the months by using month(DATE) and in the WHERE function we assigned the year value to “2015”

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

```
%sql SELECT LANDING_OUTCOME as "Landing Outcome", COUNT(LANDING_OUTCOME) AS "Total Count" FROM SPACEX \
WHERE DATE BETWEEN '2010-06-04' AND '2017-03-20' \
GROUP BY LANDING_OUTCOME \
ORDER BY COUNT(LANDING_OUTCOME) DESC ;
```

Landing Outcome	Total Count
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

By using “ORDER” we can order the values in descending order, and with “COUNT” we can count all numbers as we did previously

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 a large, brightly lit urban area is visible. In the upper right, there are greenish-yellow bands of light, likely the Aurora Borealis or Australis. The overall atmosphere is dark and mysterious.

Section 4

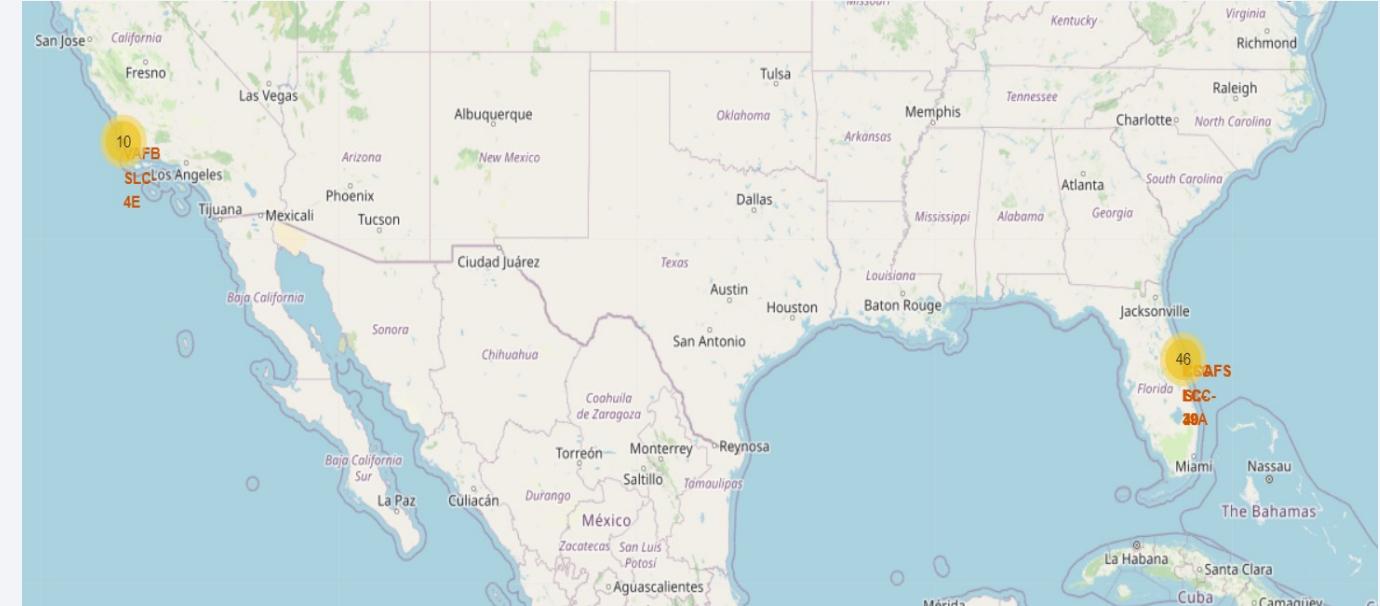
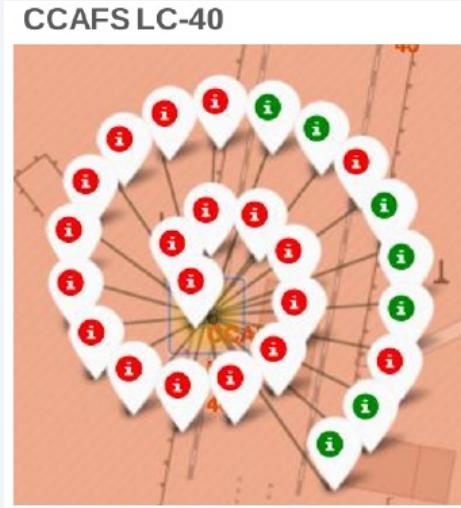
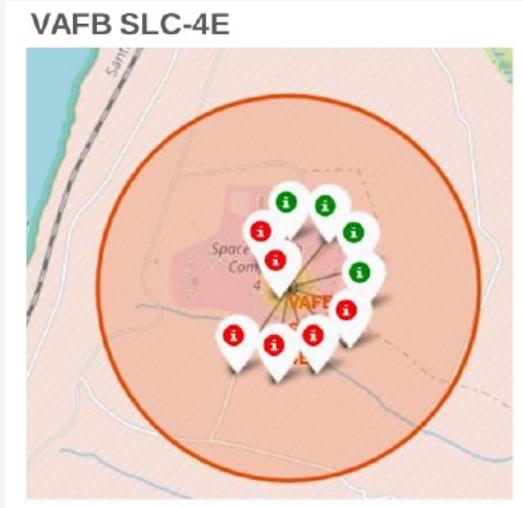
Launch Sites Proximities Analysis

All Launch Sites' Location Markers



All the launches are near USA, Florida, and California

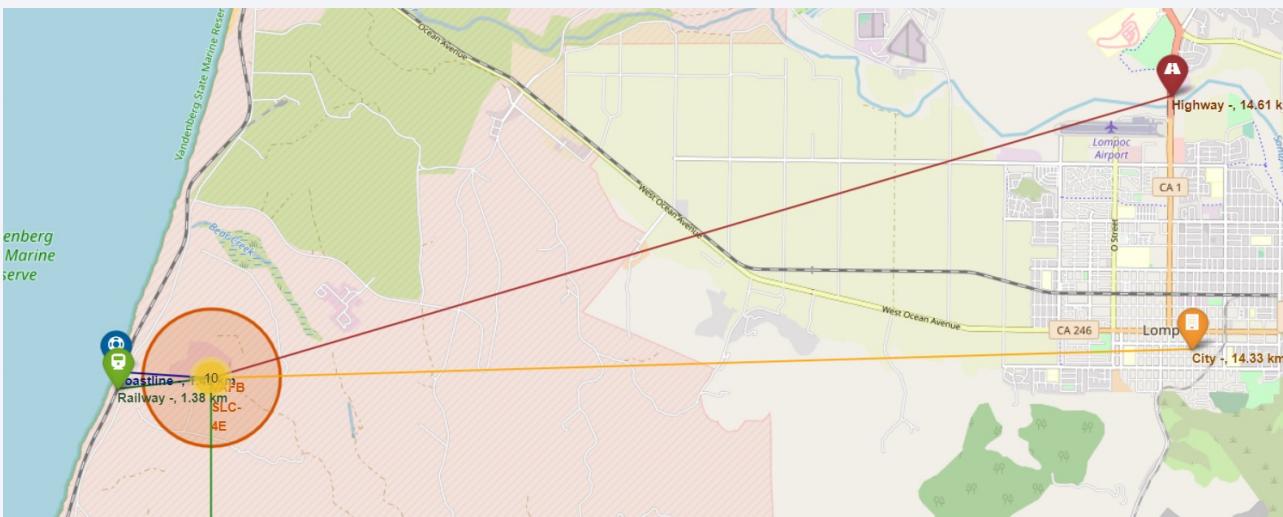
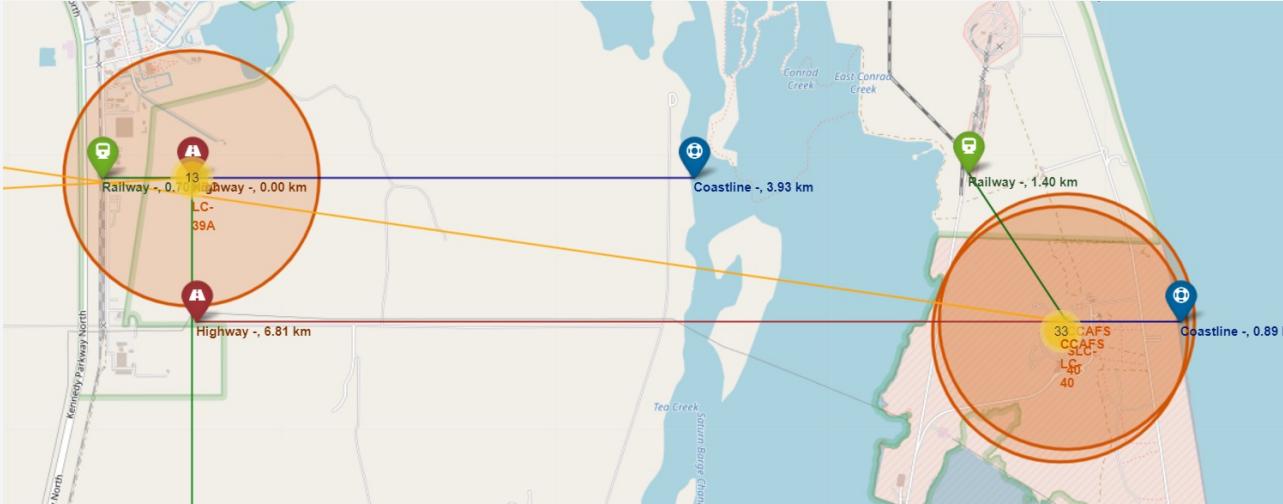
Color-labeled Launch Outcomes



Green means successful

Red means Failure

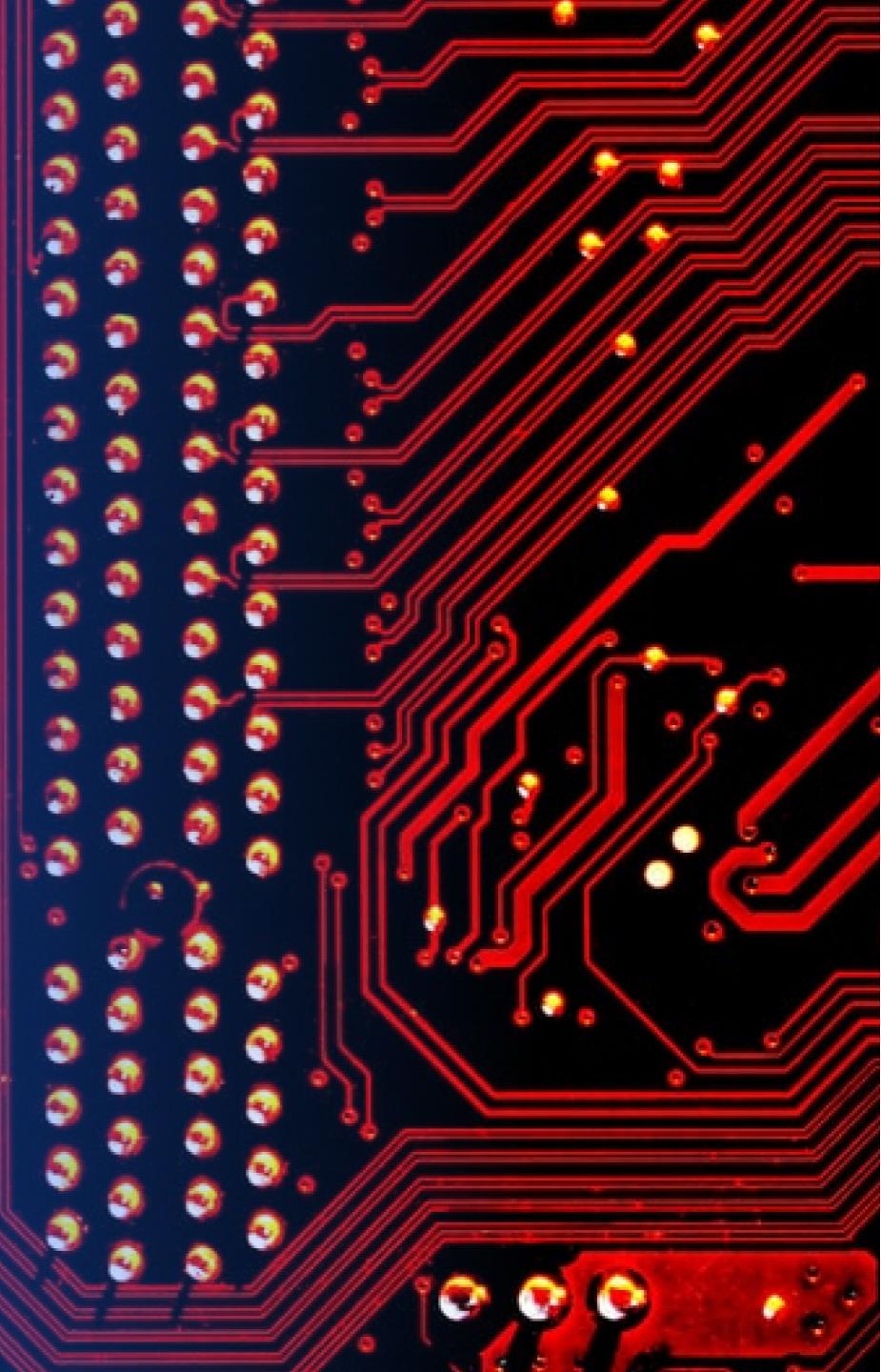
Launch Sites to its Proximities



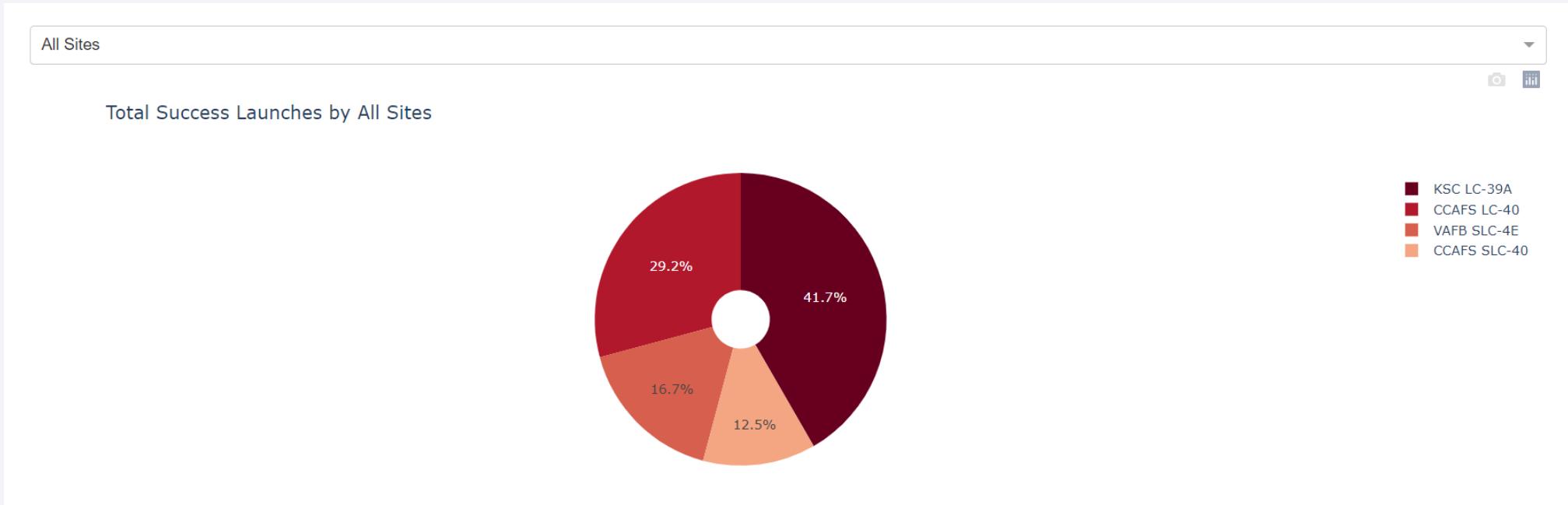
All distances from launch sites to its proximities, they weren't far from railway tracks.

Section 5

Build a Dashboard with Plotly Dash



Launch Success Count

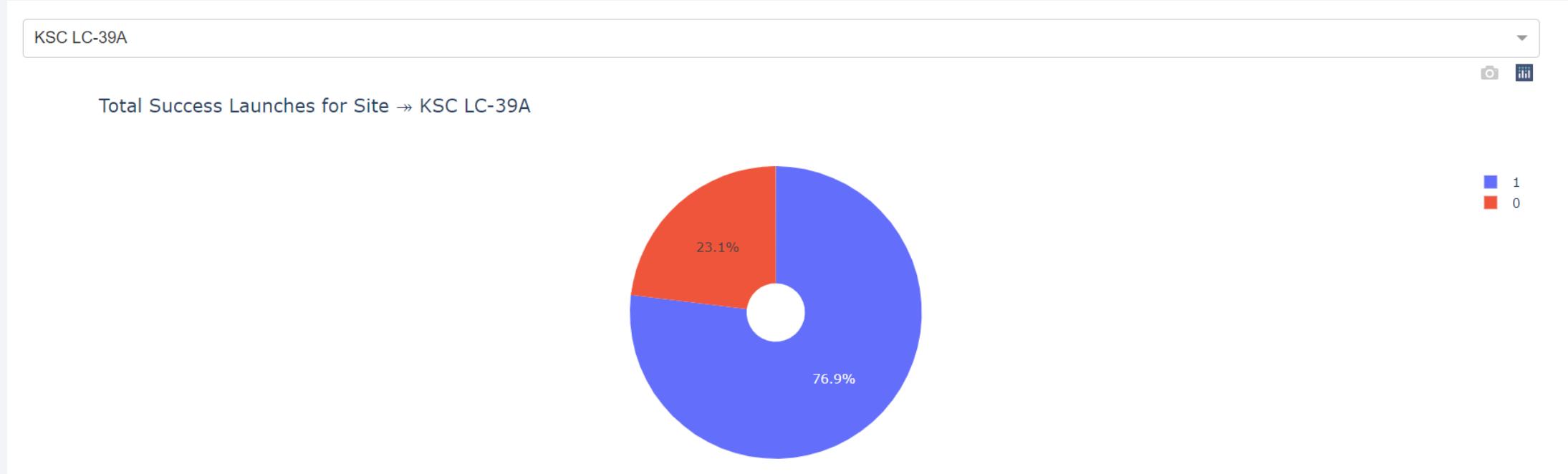


KSC LC-39A has the highest success score with 41.7%

CCAFS LC-40 comes next with 29.2%

Finally, VAFB SLC-4E and CCAFS SLC-40 with 16.7% and 12.5% respectively

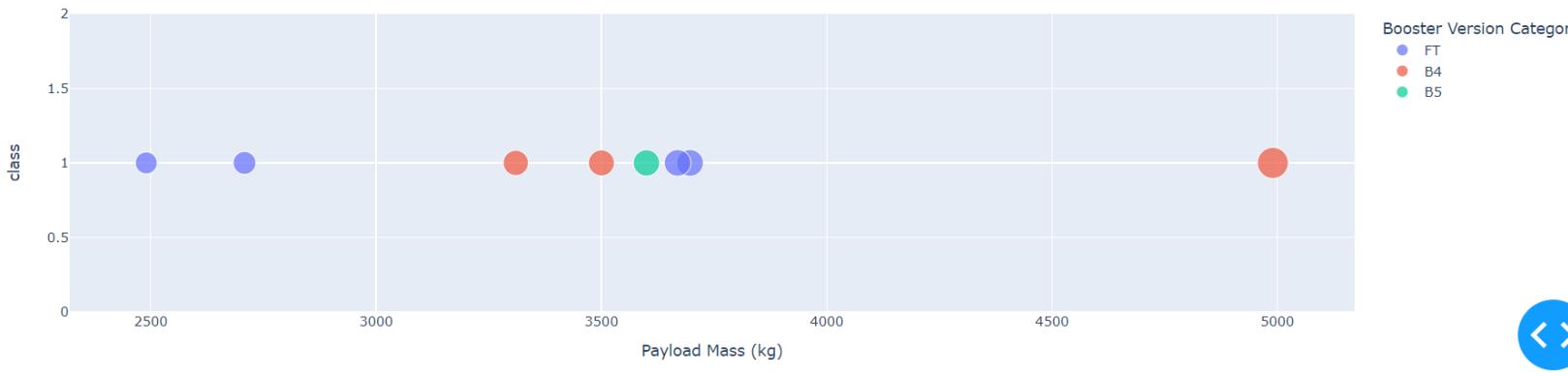
Launch Site with Highest Score



KSC LC-39A has the highest score with 76.9% with payload range of 2000 kg – 10000 kg, and FT booster version has the highest score

Payload vs. Launch Outcome

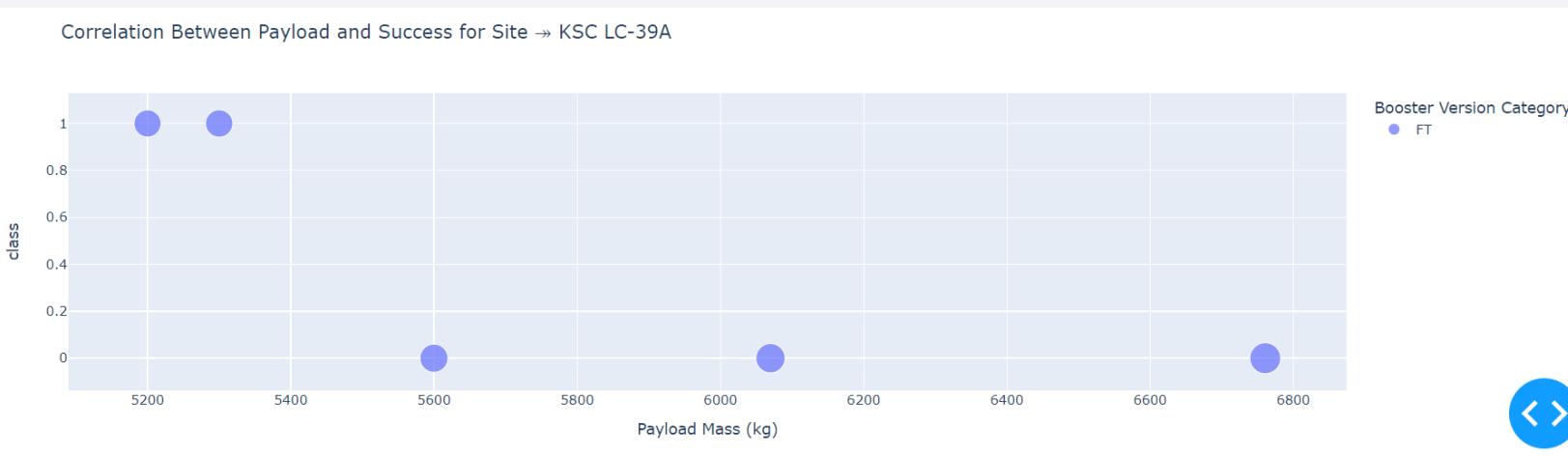
Correlation Between Payload and Success for Site → KSC LC-39A



**Payload 0 kg – 5000 kg
(first half)**



Correlation Between Payload and Success for Site → KSC LC-39A



**Payload 6000 kg – 10000 kg
(second half)**

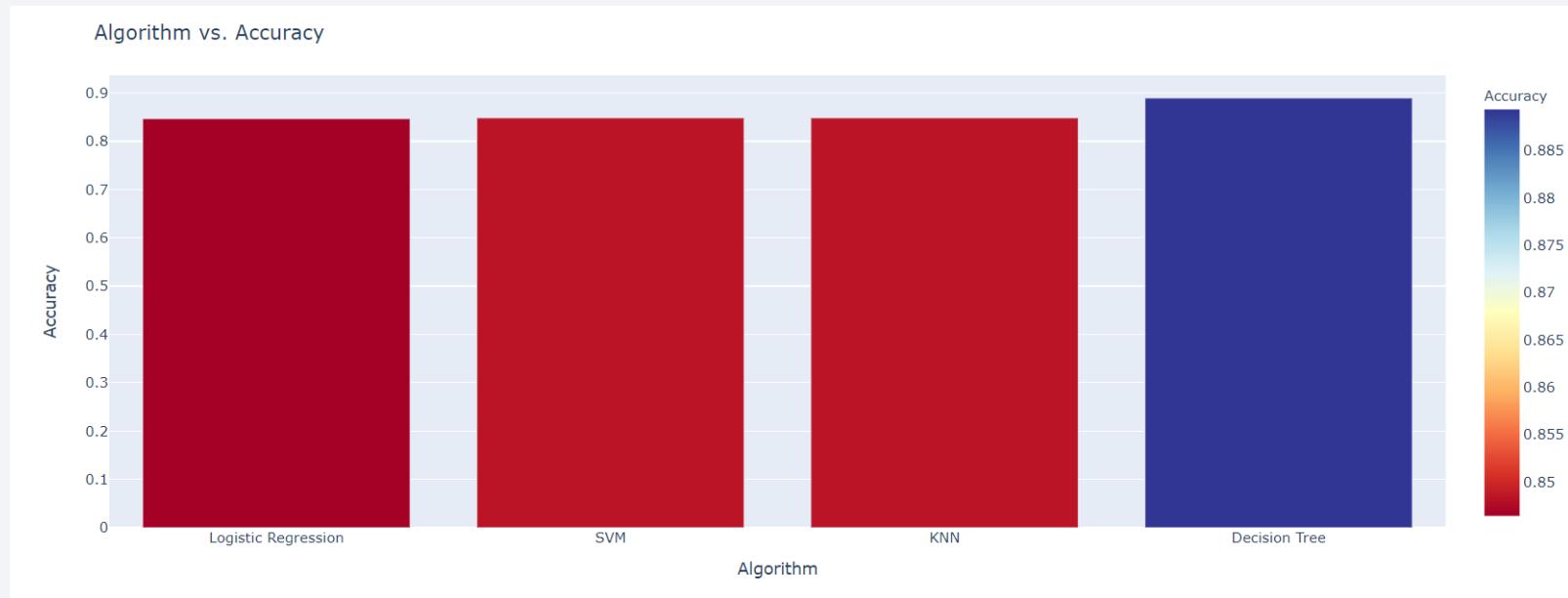


The background of the slide features a dynamic, abstract design. It consists of several thick, curved lines that transition from a bright yellow at the top right to a deep blue at the bottom left. These lines create a sense of motion and depth, resembling a tunnel or a stylized landscape. The overall effect is modern and professional.

Section 6

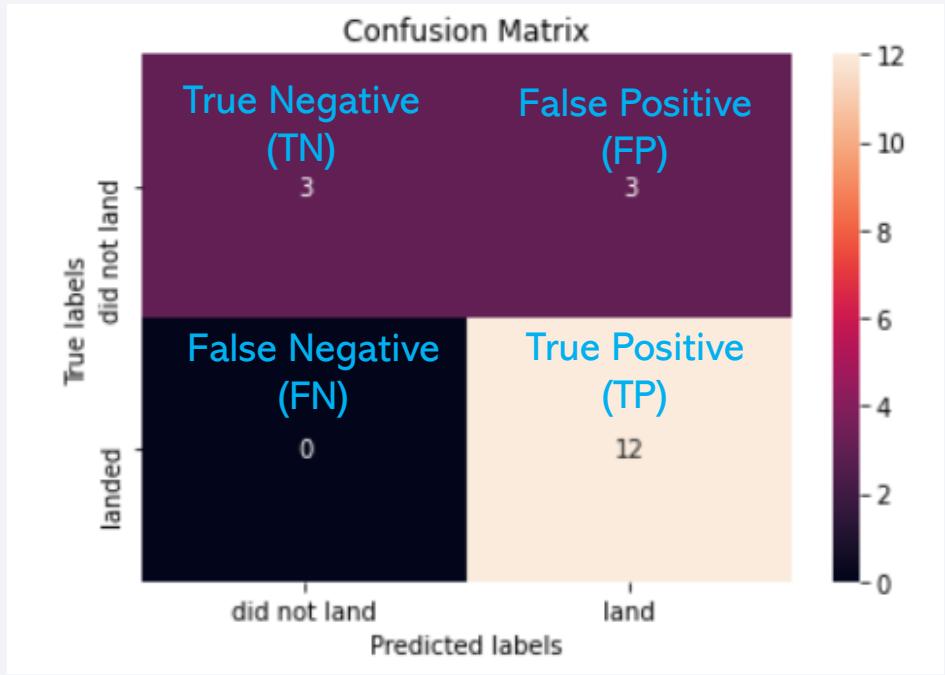
Predictive Analysis (Classification)

Classification Accuracy



Decision Tree has the highest accuracy with almost 0.89, then comes the remaining models with almost same accuracy of 0.84

Confusion Matrix



Sensitivity = 1.00, formula: $TPR = TP / (TP + FN)$

Specificity = 0.50, formula: $SPC = TN / (FP + TN)$

Precision = 0.80, formula: $PPV = TP / (TP + FP)$

Accuracy = 0.83, formula: $ACC = (TP + TN) / (P + N)$

F1 Score = 0.89, formula: $F1 = 2TP / (2TP + FP + FN)$

False Positive Rate = 0.50, formula: $FPR = FP / (FP + TN)$

False Discovery Rate = 0.20, formula: $FDR = FP / (FP + TP)$

Conclusions

- We found the site with highest score which was KSC LC-39A
- The payload of 0 kg to 5000 kg was more diverse than 6000 kg to 10000 kg
- Decision Tree was the optimal model with accuracy of almost 0.89
- We calculated the launch sites distance to its proximities

Appendix

All codes can be found on my GitHub below;

GitHub repo: <https://github.com/aynortey/IBM-Data-Science-Capstone-Project>

Thank you!

