



IBM Developer
SKILLS NETWORK

Winning Space Race with Data Science

Luciana Burdman, PhD
January 2022



Outline



Executive Summary



Introduction



Methodology



Results



Conclusion



Appendix



Executive Summary

- Summary of methodologies
 - Data Collection with Web Scraping
 - Data Wrangling
 - EDA with Data Visualization
 - EDA with SQL
 - Create Interactive Visualization with Folium
 - Create Dashboard with Plotly Dash
- Summary of all results
 - Data Analysis on the results
 - Machine Learning prediction

Introduction

- Project background and context

Commercial space age is here, companies are making space travel affordable for everyone. Perhaps the most successful companies is SpaceX. SpaceX's accomplishments include: supplying cargo to the International Space Station. Starlink, a satellite internet constellation providing internet access.

The reason why SpaceX become one of the top in this industry is because SpaceX rocket launches are relatively inexpensive. SpaceX advertised Falcon 9 rocket launches on its website at a cost of 62 million dollars; other providers cost upward of 165 million dollar each and the savings is because SpaceX can reuse the first stage of their rocket.

In this project, we would like to find out the price for each rocket launch for competitor SpaceY, competitor of SpaceX. Since the rocket cost highly depends on the successful rate of the first stage of rocket, we would like to study parameter of the First stage of rocket for successful launch.

- Project Objective - Problems you want to find answers

As a Data Scientist for SpaceY, we would like to find out how the parameters in the rocket affect on the successful landing rate of Falcon 9. Can we find out the best parameters in the first stage to ensure the successful landing rate of our First Stage rocket so that we can able to determine the best cost for the Rocket Launch.

A photograph of a large industrial building, likely a port facility, with numerous shipping containers stacked high outside. The containers are in various colors, including blue, green, red, and yellow. The building has a complex steel frame and glass windows. In the foreground, there's some foliage and a paved area.

Section 1

Methodology

Methodology

Executive Summary

Data collection methodology:

- SpaceX REST API & Web Scaping from Falcon 9
- Falcon Heavy Launches Records from Wikipedia using Beautiful Soap.

Perform data wrangling

- To convert those outcomes into Training Labels with 1 (Booster successful) (Unsuccessful)

Perform exploratory data analysis (EDA) using visualization and SQL

Perform interactive visual analytics using Folium and Plotly Dash

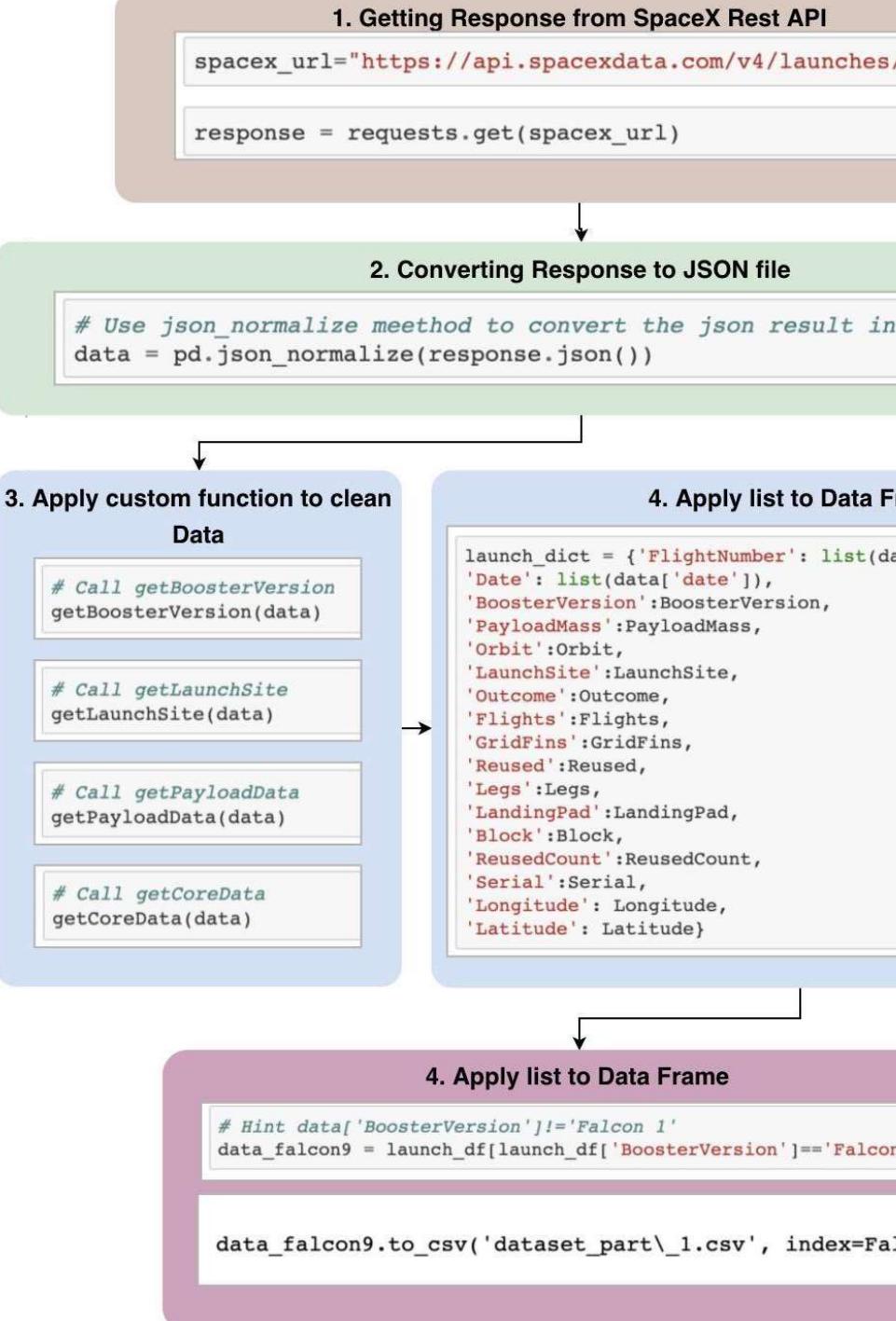
Perform predictive analysis using classification models

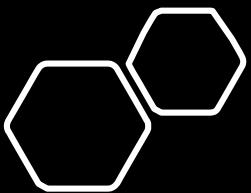
- Obtain the best Hyperparameters for SVM, Classification Trees and Logis

Data Collection

- The Data collection process includes a combination of API request from SpaceX REST API
- The API provide us information such as: Information of rocket, payload delivered. Launch specification, Landing specification, Landing outcome, location and etc.
- Using BeautifulSoup for web scraping on Wikipedia (Falcon 9 Launch data information)

Data Collection - SpaceX API





Data Collection - Web Scraping

1. Getting Response from HTML

```
static_url = "https://en.wikipedia.org/w/index.php?title=List_of_Falcon_9_and_Falcon_9_2&oldid=911911114"
```

```
# Use BeautifulSoup() to create a BeautifulSoup object from a response to static_url
soup = BeautifulSoup(html_data, 'html5lib')
```

2. Creating BeautifulSoup Object

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

3. Getting Columns Names

```
for row in first_launch_table.find_all('th'):
    name = extract_column_from_header(row)
    if(name != None and len(name) > 0):
        column_names.append(name)
```

6. Final DataFrame

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

```
extracted_row = 0
#Extract each table
for table_number,table in enumerate(soup.find_all('table')):
    # get table row
    for rows in table.find_all("tr"):
        #check to see if first table heading is as number or string
        if rows.th:
            if rows.th.string:
                flight_number=rows.th.string.strip()
                flag=flight_number.isdigit()
```

5. Appending Data

```
4.
```

```
# Remove all the tables
del launch_1
del launch_2
del launch_3
del launch_4
del launch_5
del launch_6
del launch_7
del launch_8
del launch_9
# Added new table
launch_10
launch_11
launch_12
launch_13
launch_14
```

Data Wrangling

In the data set, there are several different cases where the booster did not land successfully. Sometimes a landing was attempted but failed due to an accident; for example,

- True Ocean = mission outcome was successfully landed to a specific region of the ocean
- False Ocean = mission outcome was unsuccessfully landed to a specific region of the ocean.
- True RTLS = mission outcome was successfully landed to a ground pad
- False RTLS = mission outcome was unsuccessfully landed to a ground pad.
- True ASDS = mission outcome was successfully landed on a drone ship
- False ASDS = mission outcome was unsuccessfully landed on a drone ship.

We will mainly convert those outcomes into Training Labels with 1 means the booster successfully landed 0 means it was unsuccessful.

Perform Exploratory Data Analysis

Calculate the Number of Launch at each site

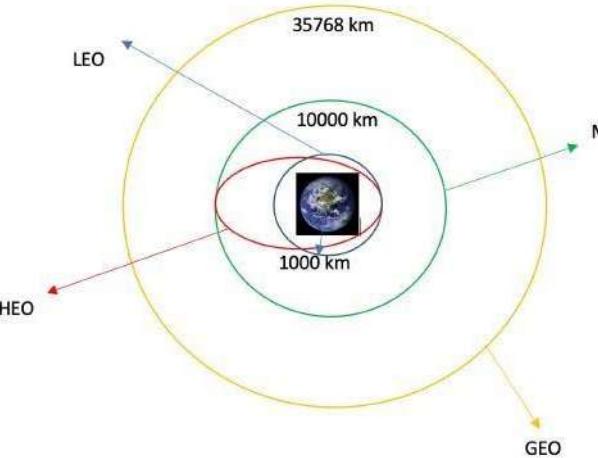
```
# Apply value_counts() on column LaunchSite  
df.value_counts(df['LaunchSite'])
```

Calculate the Number of Occurrence of mission outcome per orbit type

```
landing_outcomes = df['Outcome'].value_counts()  
landing_outcomes
```

Work out success rate for every landing

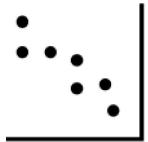
```
df['Class'].mean()  
0.6666666666666666
```



EDA with Data Visualization

Scatter Graph

1.0 Flight Number VS Launch Site



2.0 Payload VS Launch Site



3.0 Flight Number VS Orbit Type



4.0 Payload VS Orbit Type



Scatter Plot show how much one variables is affected by another. Using Scatter plot, we can check their correlation between 2 variables.

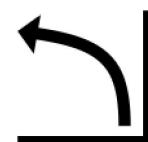
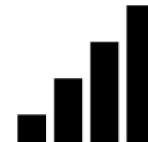
Bar Chart

1.0 Orbit Type VS Success Rate

Bar Chart make easy to compare dataset between multiple group at a glance



Bar Chart show big changes in data over time



Lin

1.0 Succe

Line Chart sh
trends very c
prediction abo

Performed SQL queries to gather information about the dataset.

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

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

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

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

Listing the date where the successful landing outcome in drone ship was achieved.

Listing the names of the boosters which have success in ground pad and have payload mass greater than 4
but less than 6000

Listing the total number of successful and failure mission outcomes

Listing the names of the booster versions which have carried the maximum payload mass.

Listing the records which will display the month names, successful landing outcomes in ground pad ,booster
versions, launch site for the months in year 2017

Ranking the count of successful landing outcomes between the date 2010-06-04 and 2017-03-20 in descen

Build an Interactive Map with Folium

Object created and added to a folium map:

Markers that show all launch sites on map

Markers that show the success/failed launches for each site on the map

Lines that show the distances between a launch site to its proximities

By adding these objects, following geographical patterns about launch sites are found:

Are launch sites in close proximity to railways? Yes

Are launch sites in close proximity to highways? Yes

Are launch sites in close proximity to coastline? Yes

Do launch sites keep certain distance away from cities? Yes

Build a Dashboard with Plotly Dash



The dashboard application contains a pie chart and a scatter point chart.



Pie Chart

- For showing total success launches by sites
- This chart can be selected to indicate a successful landing distribution across all launch sites or to indicate the success rate of individual launch sites.



Scatter Chart

- For showing the relationship between Payload mass (kg) by different categories
- Has 2 inputs: All sites/ individual launch sites on a slider between 0 and 10000
- This chart helps determine how the launch point, payload mass, and other variables affect the outcome.

Predictive Analysis (Classification)



Building Model

Load our dataset into Numpy and Pandas
Transform Data
Split our data into training and test data sets
Check how many test samples
Decide on which type of machine learning algorithms to apply



Evaluating Model

Check accuracy for each model
Get tuned hyperparameters for each type of algorithms



Improving Model

Feature Engineering
Algorithm Tuning



Finding The best Classification Model

The model with the highest accuracy wins the best performance award.

Results



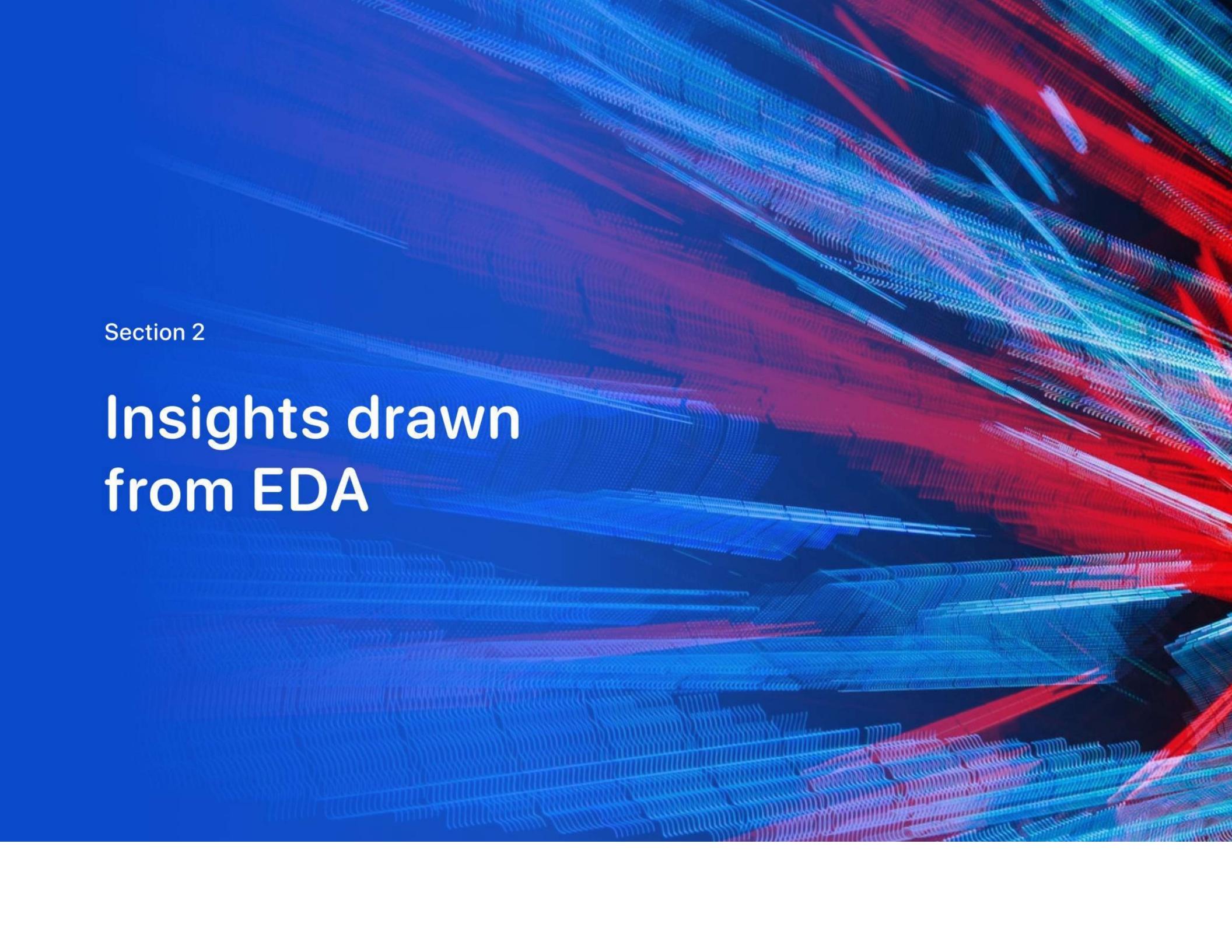
EXPLORATORY DATA
ANALYSIS RESULTS



INTERACTIVE ANALYTICS
DEMO IN SCREENSHOTS

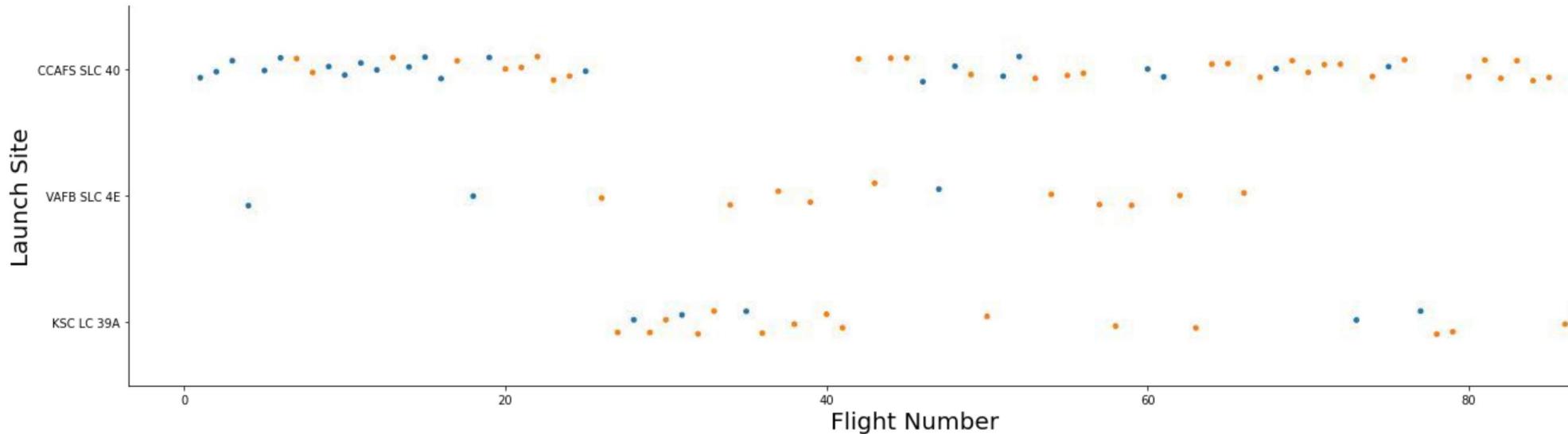


PREDICTIVE A
RESULTS

The background of the slide features a complex, abstract pattern of wavy, horizontal lines in shades of blue, red, and white. These lines create a sense of depth and motion, resembling a digital or architectural landscape. They are more concentrated in the lower half of the slide, with some lines extending towards the top right corner.

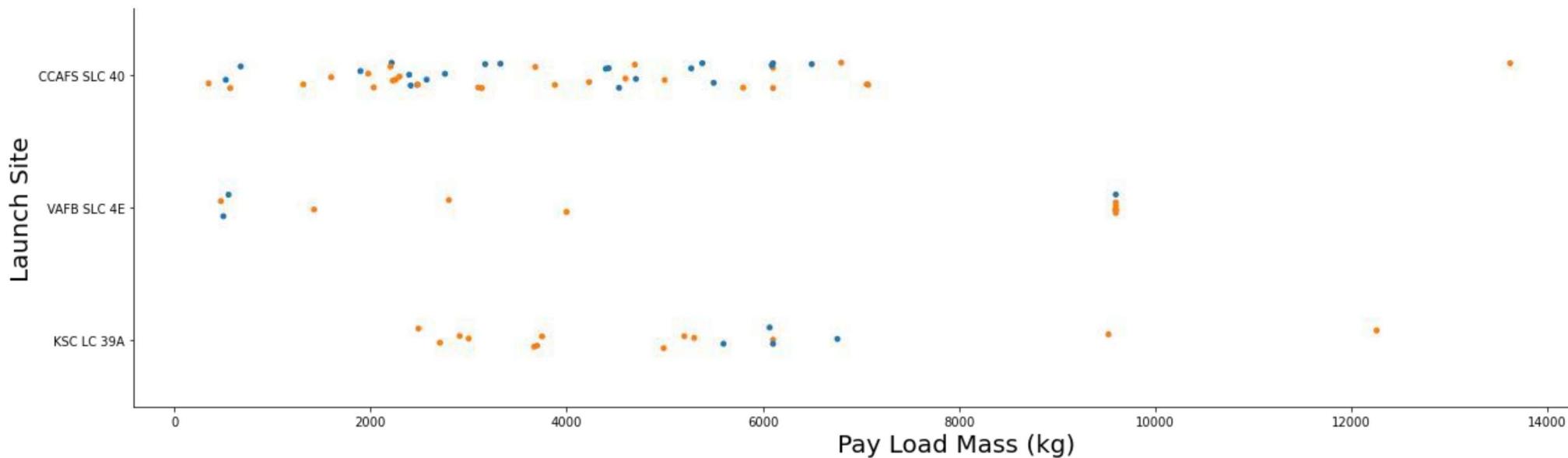
Section 2

Insights drawn from EDA



- Class 0 (Blue) represents unsuccessful launch, and Class 1 (orange) represents successful launch.
- This figure shows that the success rate increased as the number of flights increased
- As the success rate has increased considerably since the 20th flights. This point seems to be a big breakthrough

Flight Number vs. Launch Site

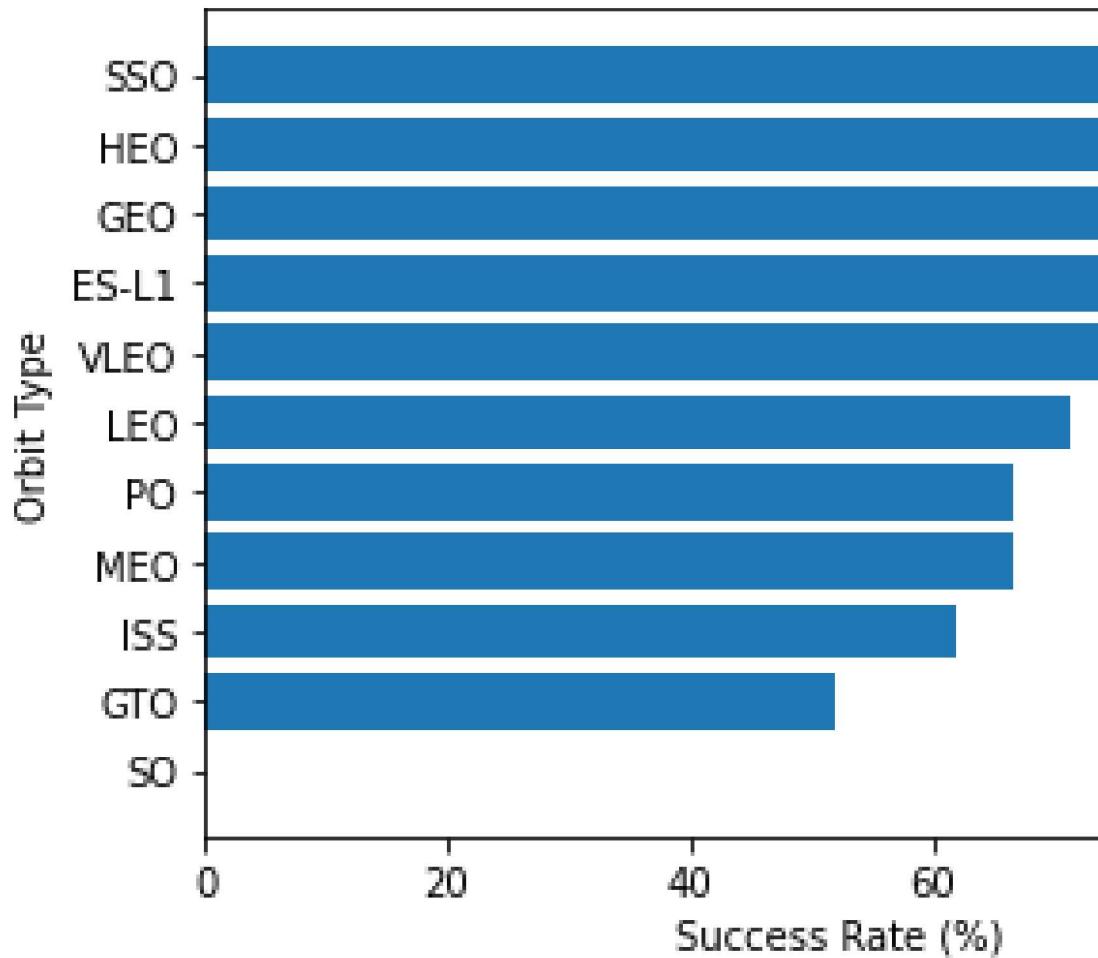


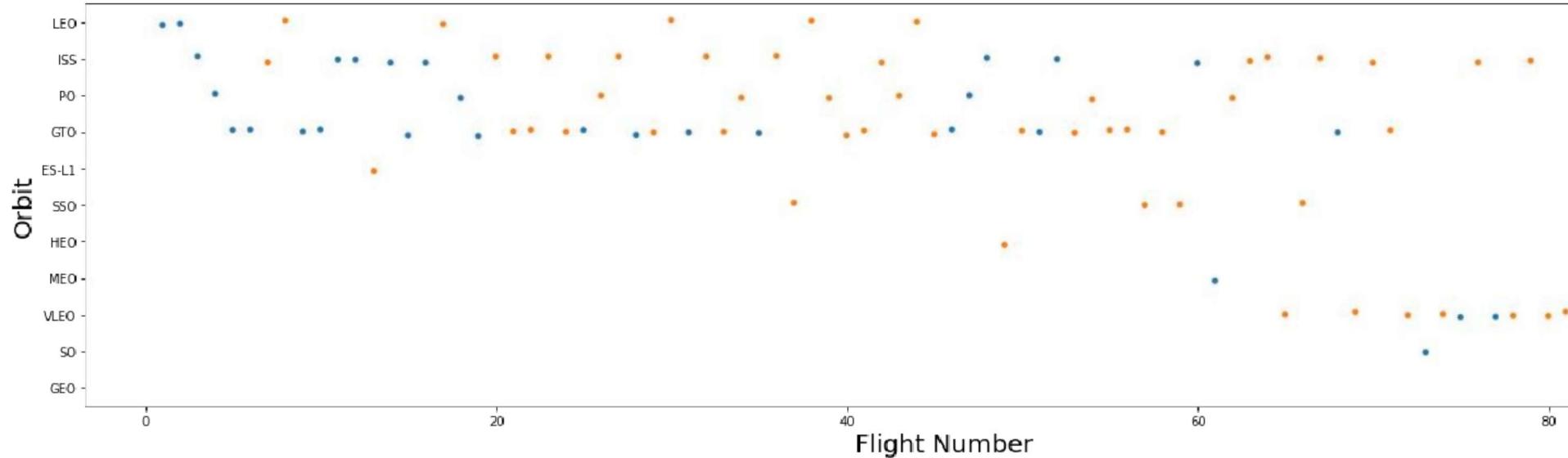
- Class 0 (Blue) represents unsuccessful launch, and Class 1 (orange) represents successful launch.
- At first glance, the larger pay load mass, the higher the rocket's success rate, but it seems difficult to make decisions based on this figure because no clear pattern can be found between successful launch Pay Load Mass

Payload vs. Launch Site

Success Rate vs. Orbit Type

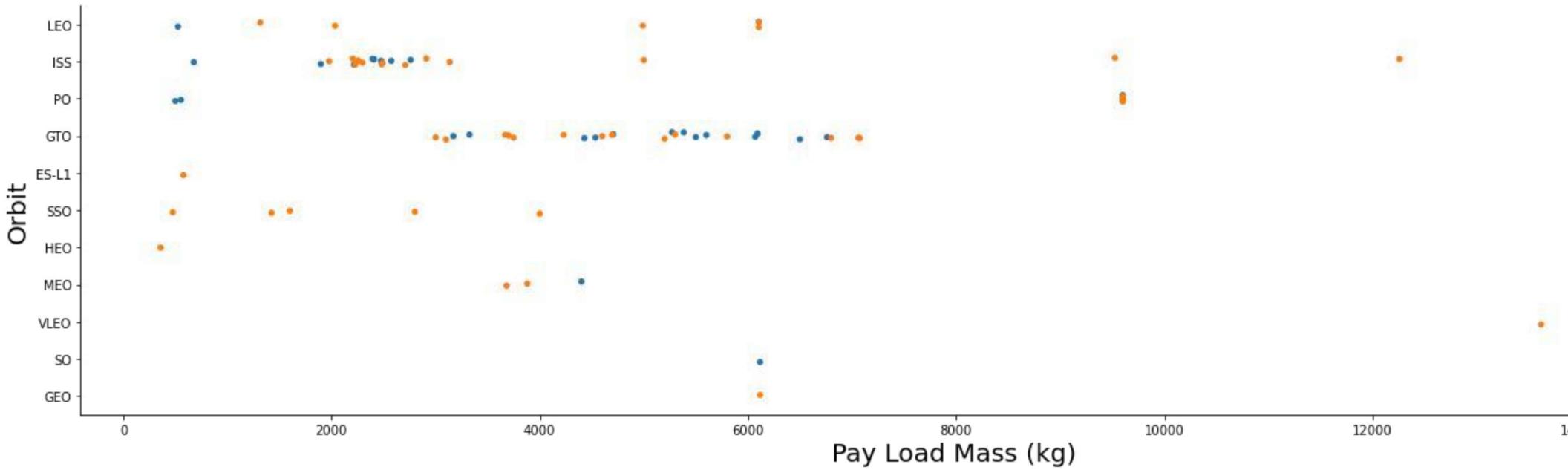
- Orbit types SSO, HEO, GEO, and ES-L1 have the highest success rates (100%)
- On the other hand, the success rate of orbit type GTO is only 50%, and it is the lowest except for type SO, which recorded failure in a single attempt.





- Class 0 (Blue) represents unsuccessful launch, and Class 1 (orange) represents successful launch.
- In most cases, the launch outcome seems to be correlated with the flight number.
- On the other hand, in GTO orbit there seems to be no relationship between flight numbers and success rate.
- SpaceX starts with LEO with a moderate success rate, and it seems that VLEO, which has a high success rate, is recent launches

Flight Number VS Orbit Type



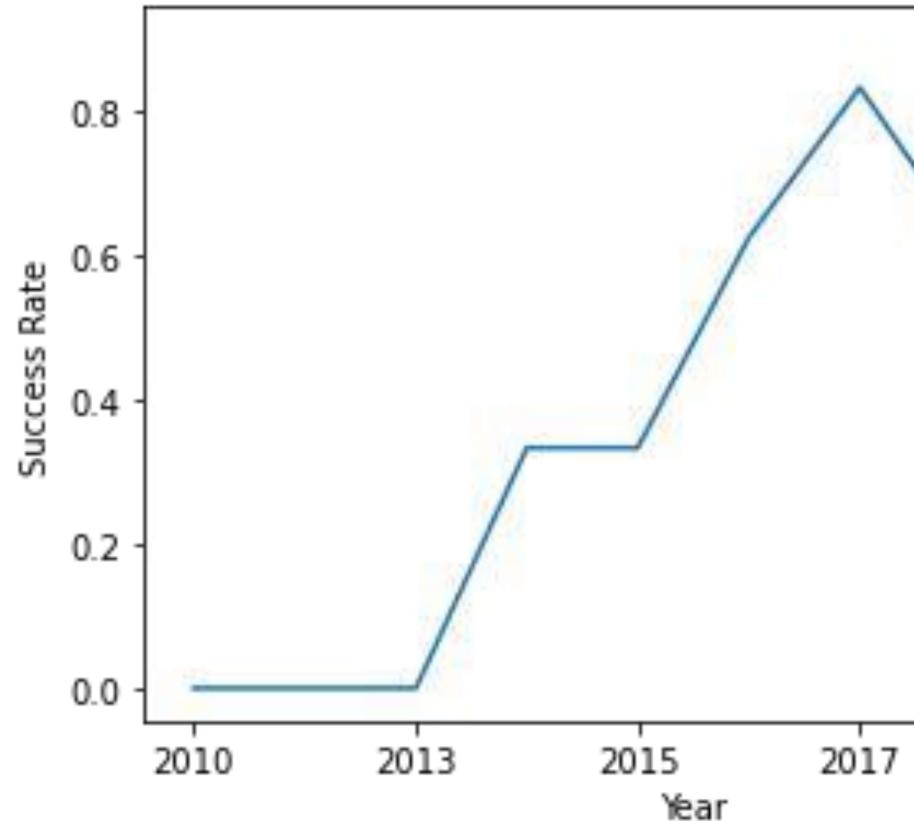
- Class 0 (Blue) represents unsuccessful launch, and Class 1 (orange) represents successful launch.
- With heavy payloads, the successful landing rate are higher for LEO and ISS.
- However, for GTO case, it is hard to distinguish between the positive landing rate and the negative landing because gathered.

Payload VS Orbit Type

Launch Success

Yearly Trend

- Since 2013, the success rate has continued to increase until 2017
- The rate decreased slightly in 2018
- Recently, it has shown a success rate of about 80%



All Launch Site Names

GITHUB

%%sql

```
SELECT DISTINCT LAUNCH_SITE  
FROM SPACEXTBL
```

launch_site

CCAFS LC-40

CCAFS SLC-40

KSC LC-39A

VAFB SLC-4E



Query Explanation



Using the word DISTINCT in the query means that it will return unique values in the Launch_Site column from SpaceX data.



There are four unique launch sites: CCAFS LC-40, CCAFS SLC-40, KSC LC-39A, VAFB SLC-4E

Launch Site Names Begin with 'CCA'

- 5 records of the SpaceX table were displayed using LIMIT 5 clause in the query
- Using the LIKE operator and the percent sign (%) together, the LAUNCH_SITE name starting with CAA will be called

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

DATE	time_utc	booster_version	launch_site	payload	payload_mass_kg	orbit	cu
2010-06-04	18:45:00	F9 v1.0 B0003	CCAFS LC-40	Dragon Spacecraft Qualification Unit	0	LEO	
2010-12-08	15:43:00	F9 v1.0 B0004	CCAFS LC-40	Dragon demo flight C1, two CubeSats, barrel of Brie cheese	0	LEO (ISS)	NASA
2012-05-22	07:44:00	F9 v1.0 B0005	CCAFS LC-40	Dragon demo flight C2	525	LEO (ISS)	NASA
2012-10-08	00:35:00	F9 v1.0 B0006	CCAFS LC-40	SpaceX CRS-1	500	LEO (ISS)	NASA
2013-03-01	15:10:00	F9 v1.0 B0007	CCAFS LC-40	SpaceX CRS-2	677	LEO (ISS)	NASA

Total Payload Mass

- Using the SUM() function to calculate the sum of column PAYLOAD_MASS_KG
- The WHERE Clause filter the dataset to only perform calculations on CUSTOMER NASA (CRS)

```
%%sql
SELECT SUM(PAYLOAD_MASS__KG_) AS TOTAL
FROM SPACEXTBL
WHERE CUSTOMER = 'NASA (CRS)'
```

total_payload_mass_kg

45596

Average Payload Mass by F9v1.1

- Using the **AVG()** function to calculate the average value of column **PAYLOAD_MASS_KG**
- The WHERE clause filters the dataset to only perform calculation on **Booster_version = F9 v1.1**

```
%%sql
SELECT AVG(PAYLOAD_MASS__KG_) AS AVG_
FROM SPACEXTBL
WHERE BOOSTER_VERSION = 'F9 v1.1'
```

avg_payload_mass_kg

2928

First Successful Ground Landing Date

- Using the **MIN()** function to find out the earliest date in column **DATE**
- The WHERE clause filters the dataset to only perform filtration on **LANDING_OUTCOME**

```
SELECT MIN(DATE) AS FIRST_SUCCESSFUL_LANDING_DATE  
FROM SPACEXTBL  
WHERE LANDING_OUTCOME = 'Success' (group by)
```

fist_successful_landing_date

2015-12-22

Successful drone ship landing with payload between 4000 and 6000

- Selecting only BOOSTER_VERSION
- The WHERE clause filters the dataset to LANDING_OUTCOME = Success (drone ship)
- The AND and BETWEEN clause specifies additional filter condition PAYLOAD_MASS_KG BETWEEN 4000 AND 6000

```
%%sql
SELECT BOOSTER_VERSION FROM SPACEXTBL
WHERE LANDING_OUTCOME = 'Success' (drone ship)
AND (PAYLOAD_MASS_KG BETWEEN 4000 AND 6000)
```

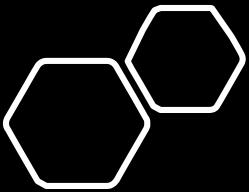
booster_version

F9 FT B1022

F9 FT B1026

F9 FT B1021.2

F9 FT B1031.2



Total Number of Successful and Failure Mission Outcomes

- Using the COUNT() function to filter the total number of columns
- Using GROUP BY function to group rows that have same values into summary rows to find the total number in each MISSION_OUTCOME
- SpaceX successfully completed nearly 99% of its mission based on the dataset

```
%%sql
SELECT MISSION_OUTCOME, COUNT(*) AS COUNT_OF_MISSIONS
FROM SPACEXTBL
GROUP BY MISSION_OUTCOME
```

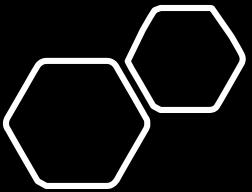
mission_outcome	total_missions
Failure (in flight)	1
Success	100
Success (payload status unclear)	1

Boosters Carried Maximum Payload

- Using a subquery, find the maximum value of the payload using MAX() function, and then filter the dataset to perform search IF PAYLOAD_MASS_KG_ is the maximum value
- From the result, F9 B5 B10xx.x boosters carried the maximum payload

```
%%sql
SELECT DISTINCT BOOSTER_VERSION, P
FROM SPACEXTBL
WHERE PAYLOAD_MASS_KG_ = (
    SELECT MAX(PAYLOAD_MASS_KG_)
    FROM SPACEXTBL);
```

booster_version	payload_mass_kg_
F9 B5 B1048.4	15600
F9 B5 B1048.5	15600
F9 B5 B1049.4	15600
F9 B5 B1049.5	15600
F9 B5 B1049.7	15600
F9 B5 B1051.3	15600
F9 B5 B1051.4	15600
F9 B5 B1051.6	15600
F9 B5 B1056.4	15600
F9 B5 B1058.3	15600
F9 B5 B1060.2	15600
F9 B5 B1060.3	15600

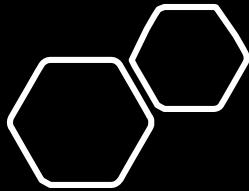


2015 Launch Records

- In the WHERE clause, filter the dataset to perform a search if Landing_Outcome is Failure (drone ship)
- Use AND operator to display a record if YEAR is 2015
- There were two landing failures on drone ships in 2015

```
%%sql
SELECT LANDING__OUTCOME, BOOSTER_VERSION, LAUNCH_DATE
FROM SPACEXTBL
WHERE LANDING__OUTCOME = 'Failure (drone ship)'
```

landing__outcome	booster_version	launched
Failure (drone ship)	F9 v1.1 B1012	2015-01-01
Failure (drone ship)	F9 v1.1 B1015	2015-01-02



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

- In the WHERE clause, filter the dataset to perform search for DATE between 2010-06-04 and 2017-03-20
- Using ORDER clause to sort the records by total number of landing and DESC clause to sort the records in descending order.

```
%%sql
SELECT LANDING_OUTCOME, COUNT(LANDING_OUTCOME)
FROM SPACEXTBL
WHERE DATE BETWEEN '2010-06-04' AND '2017-03-20'
GROUP BY LANDING_OUTCOME
ORDER BY total_number DESC
```

landing_outcome	total_number
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

The background of the slide is a nighttime satellite photograph of Earth. The dark blue oceans are visible, along with the glowing lights of numerous cities and towns across the continents. The atmosphere appears as a thin, hazy layer above the surface.

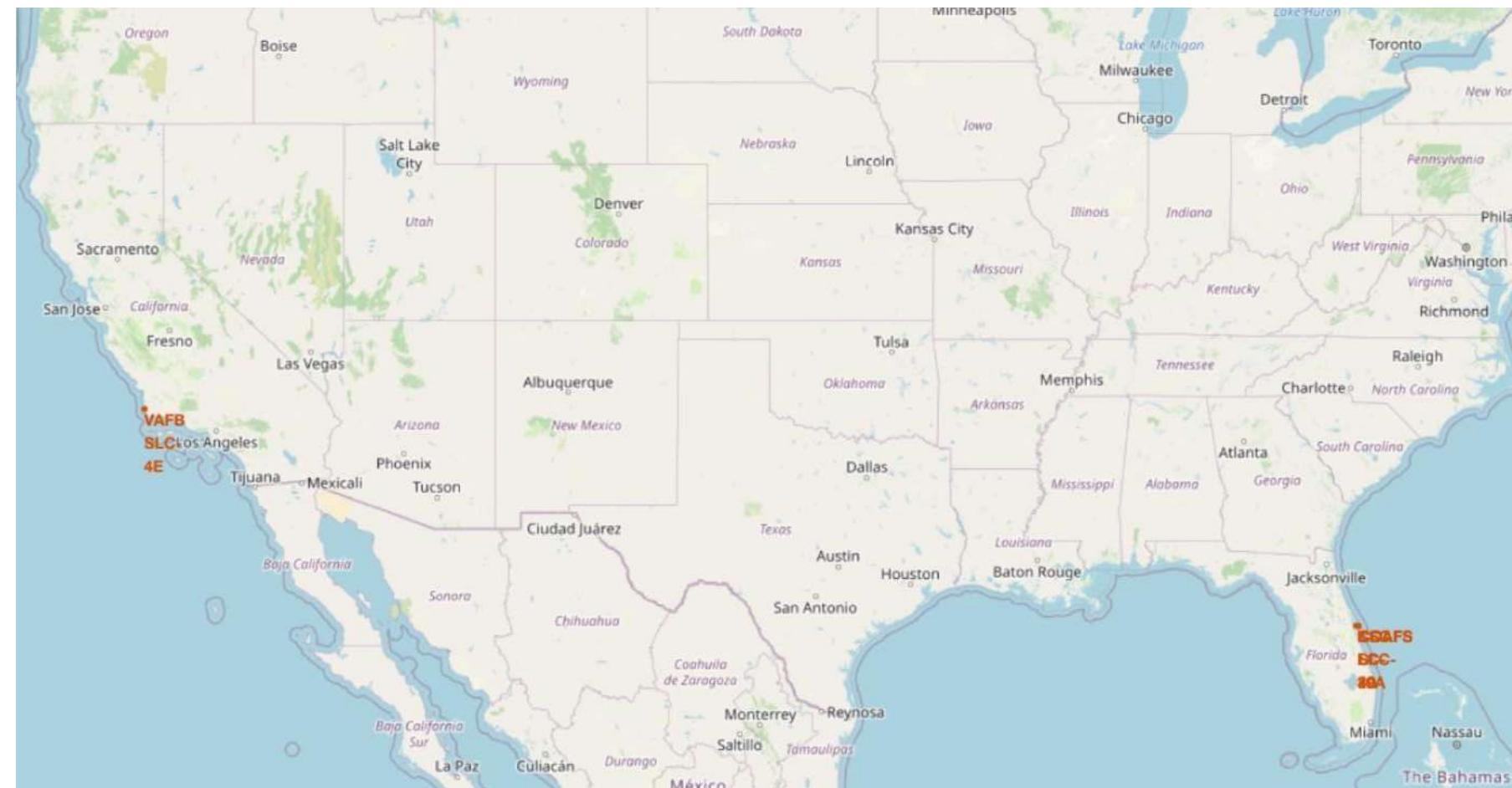
Section 4

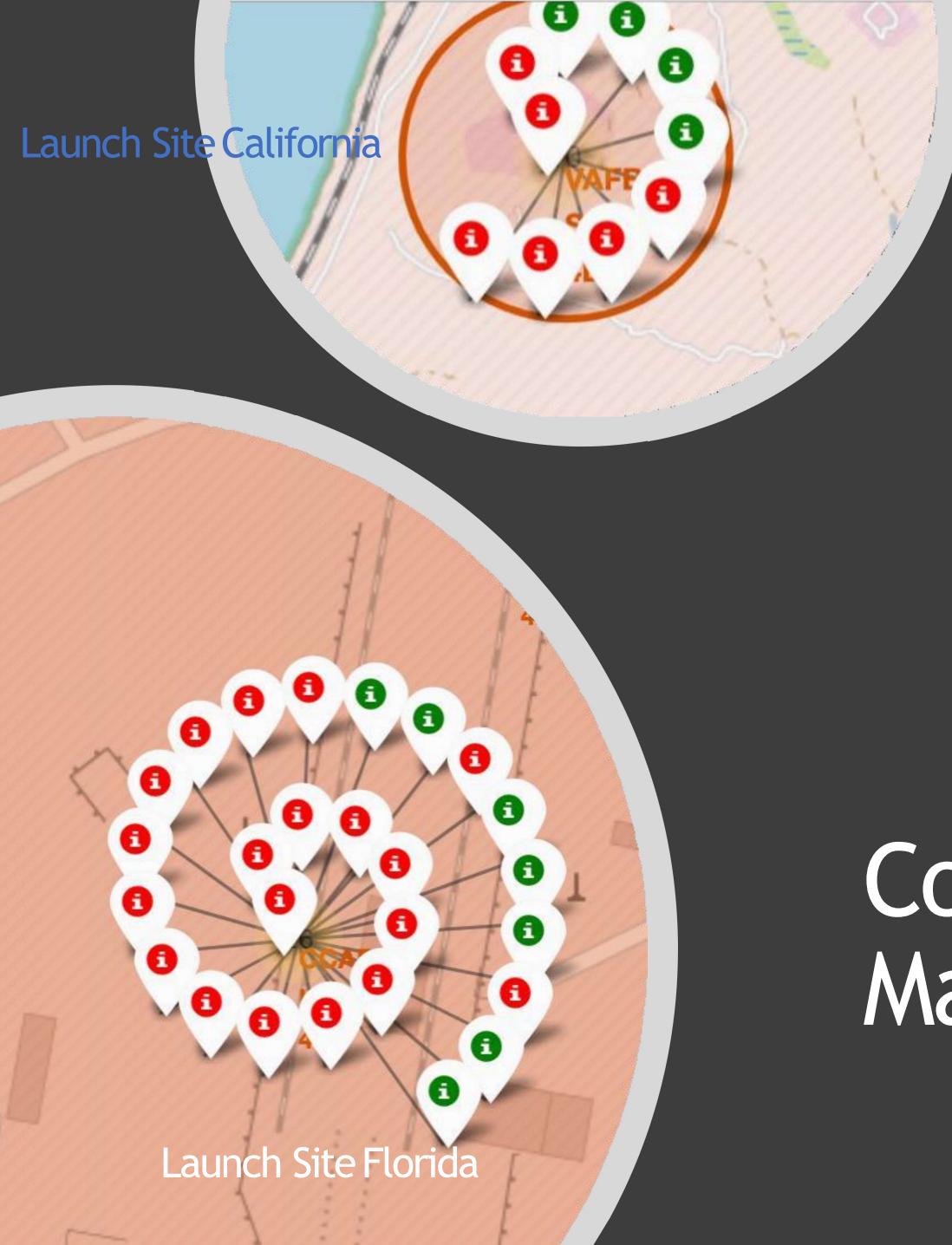
Launch Sites Proximities Analysis

All Launch Site Location

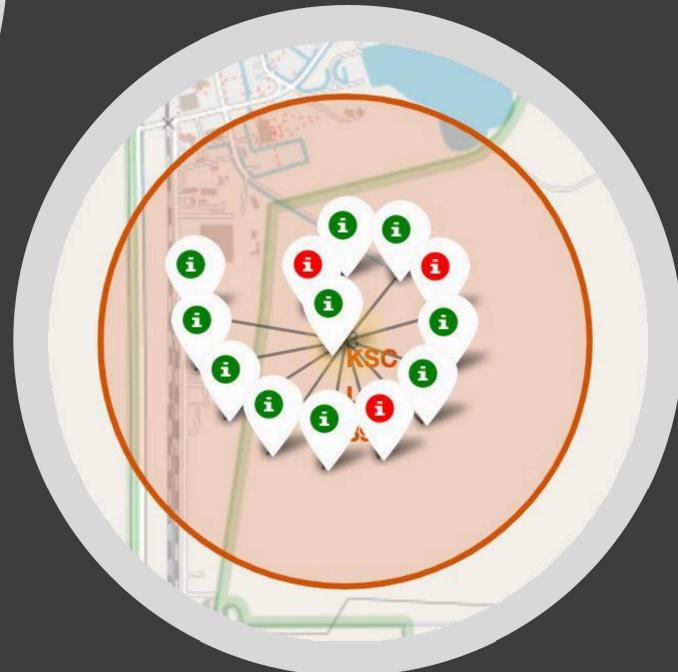
Most of the site are in

Florida and



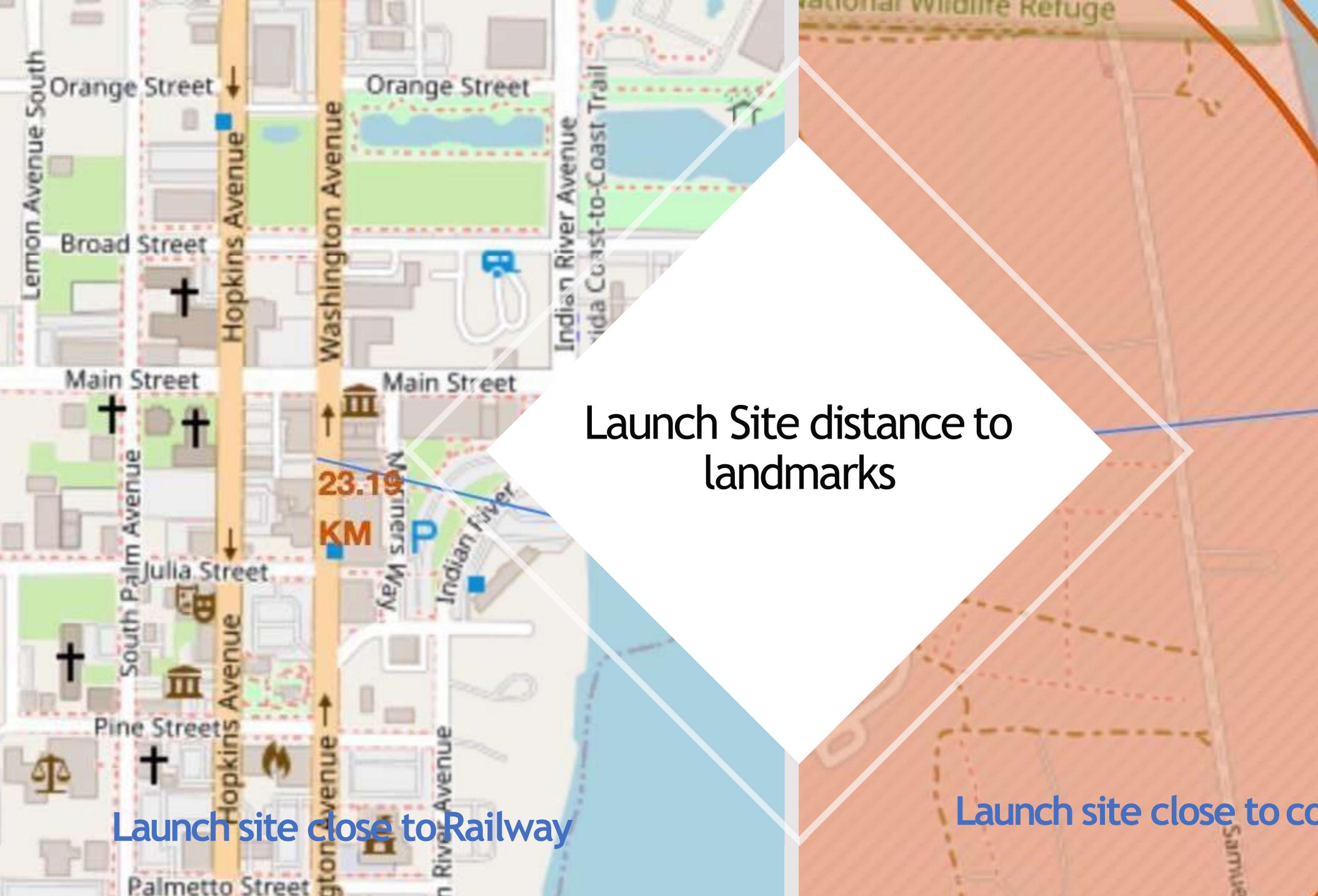


Launch Site Florida



Colour Labelled Markers





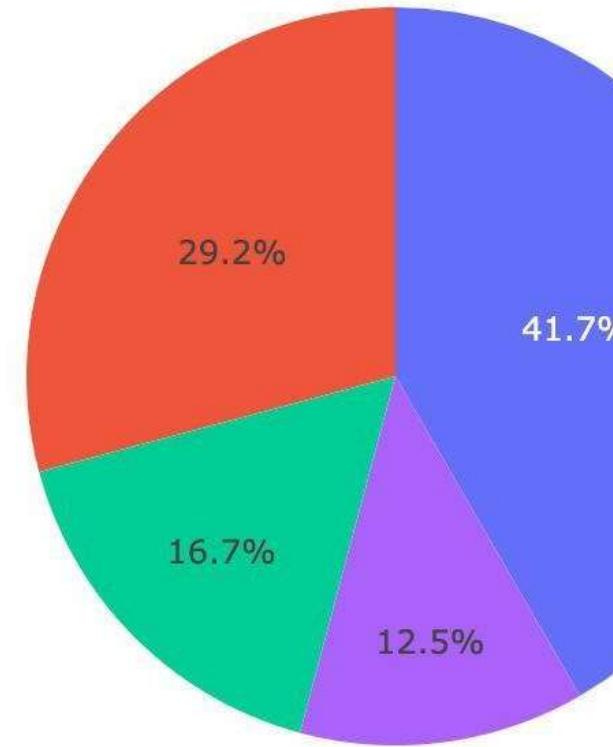
Launch Site distance to landmarks

Section 5

Build a Dashboard with Plotly Dash

Total Success Launches by all Sites

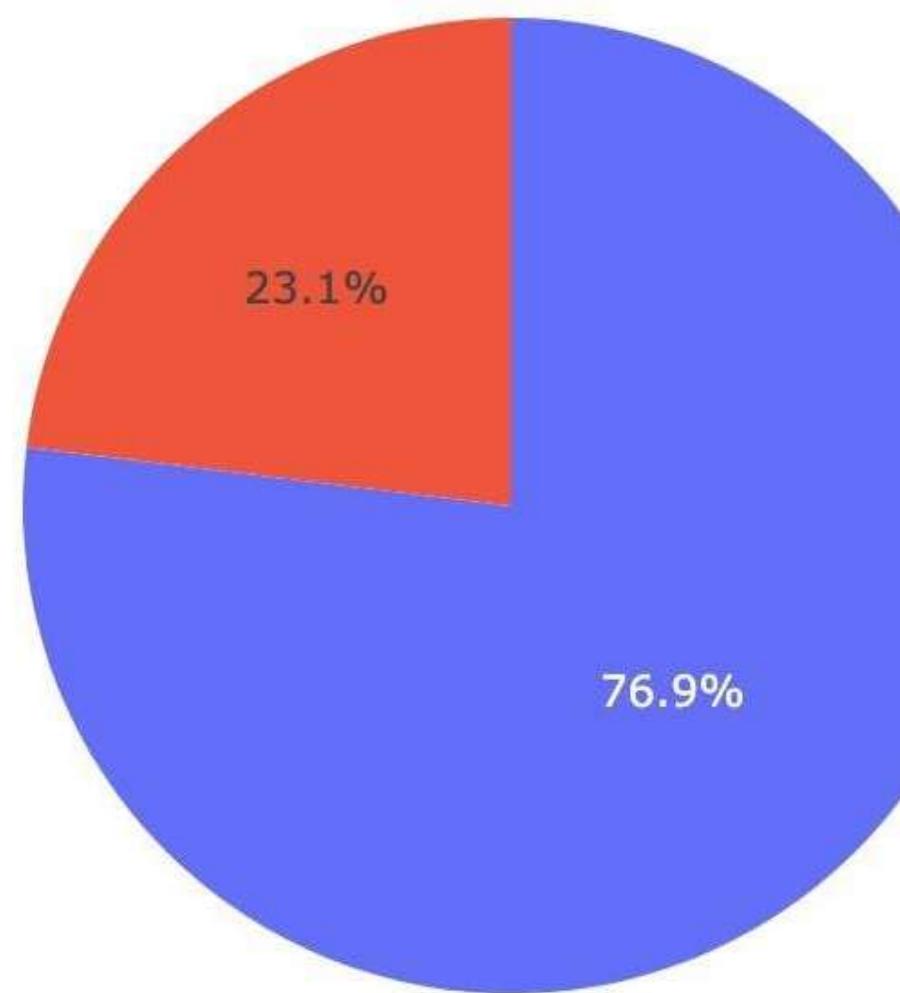
- KSLC - 39A records the most launch success among all sites.
- VAFB SLC-4E has the lowest success launch



■ KSC LC-39A
■ CCAFS LC-4
■ VAFB SLC-4
■ CCAFS SLC-

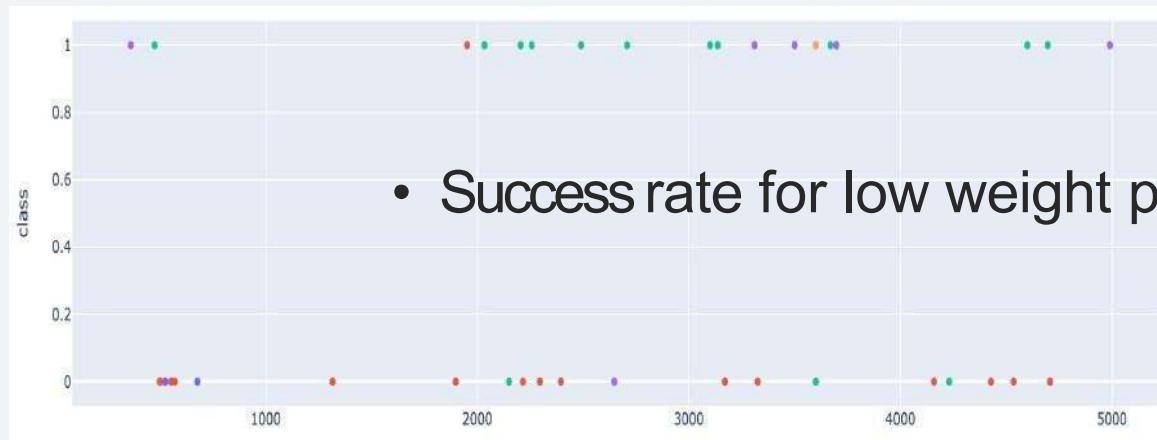
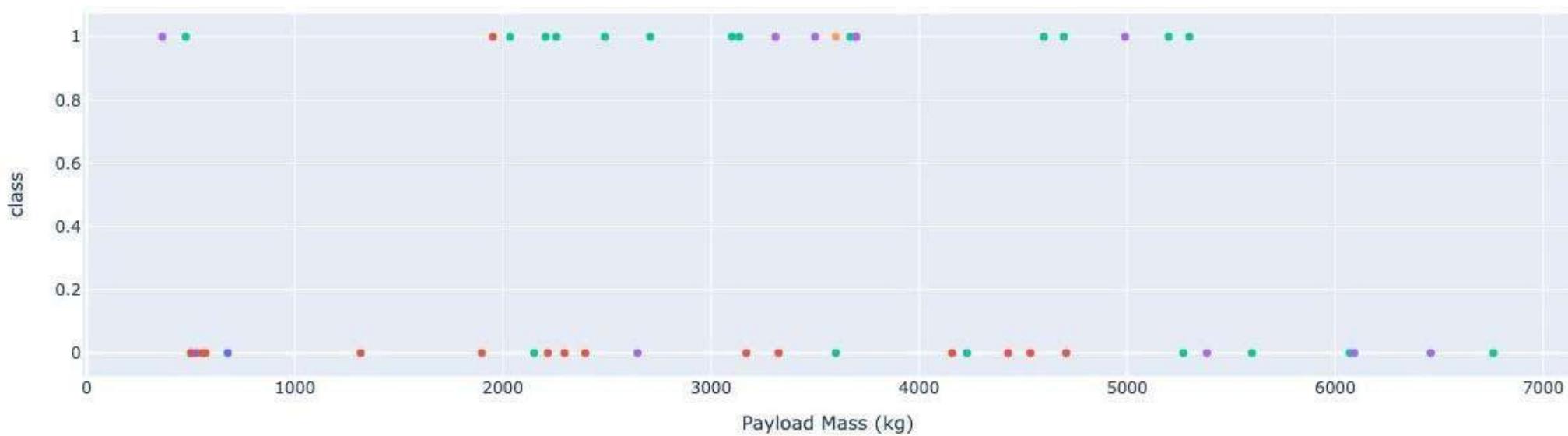
Launch Site with Highest launch Success Ratio

- KSC-LC-39A achieved a 76.9% success rate with total of 13 landing

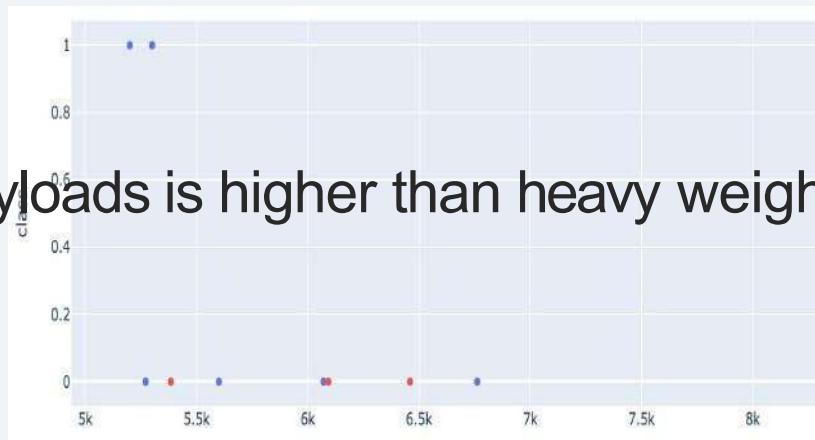


Payload VS Launch Outcome Scatter Plot for all Sites

Correlation between Payload and Success for all Sites



- Success rate for low weight payloads is higher than heavy weight

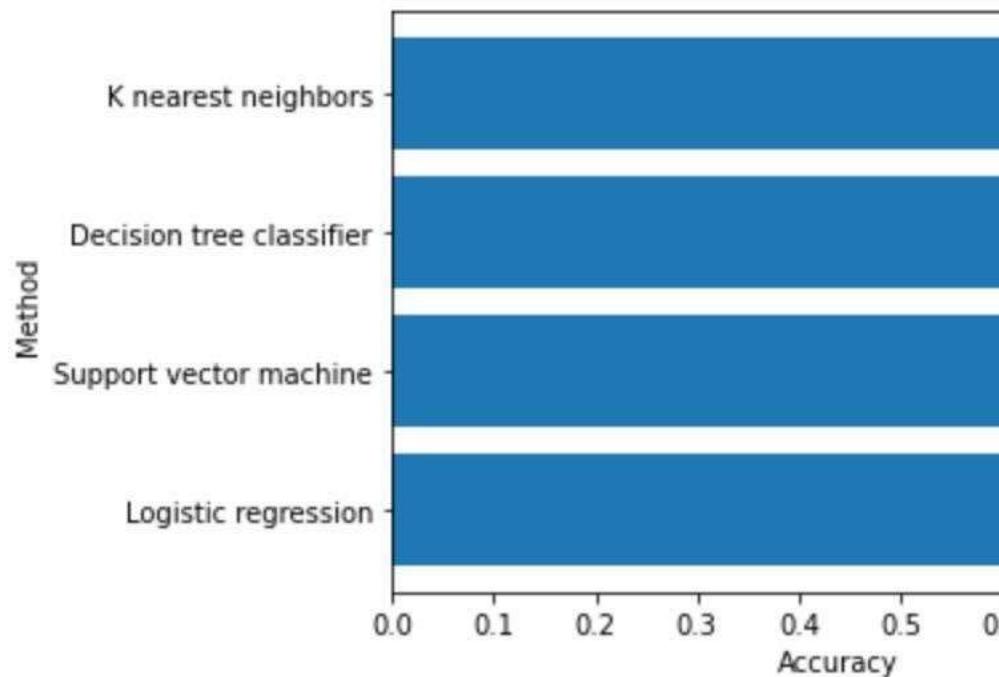


Section 6

Predictive Analysis (Classification)

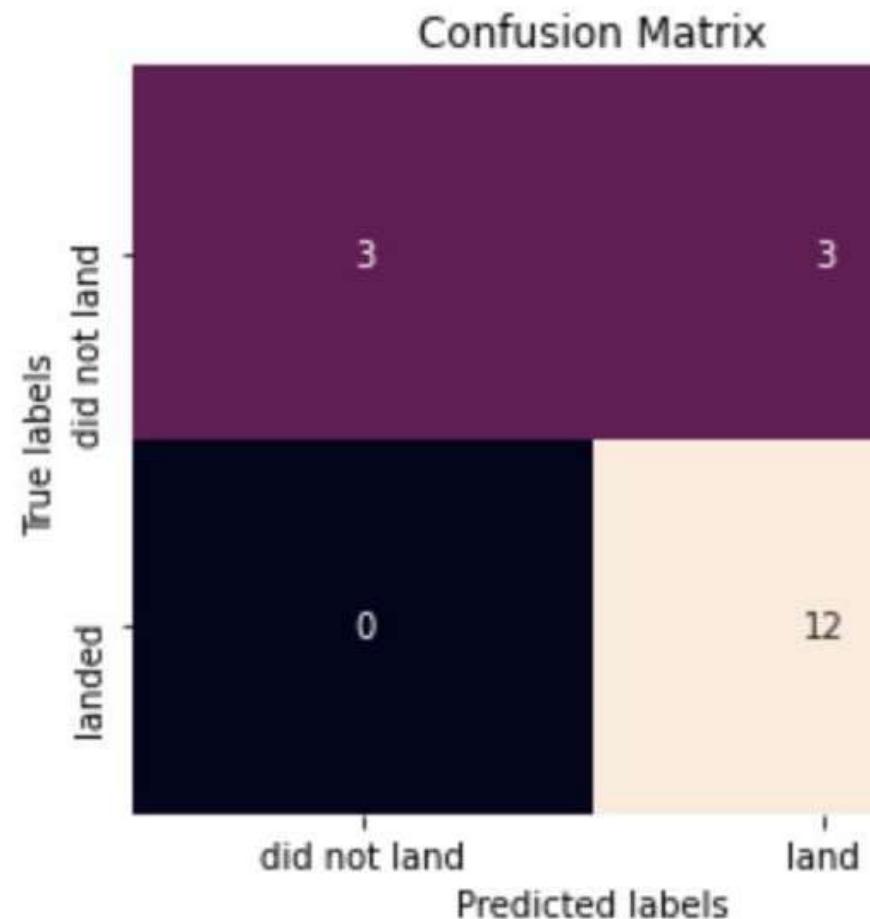
Classification Accuracy

- In the test set, the accuracy of all models was virtually the same at 83.33%
- More data is needed to improve the model



Confusion Matrix

- The confusion matrix is the same for all models
- The models predicted 12 successful landing when the true label was successful.
- The model predict 3 failed landing while the actual label was not successful.
- Overall, the model predict quite good at successful landings.





Orbital types SSO, HEO, GEO, and ES-L1 have the highest success rate (100%).



KSLC-39A has the highest number of launch successes and the highest success rate among all sites.



Low weighted payloads perform better than the heavier payloads



In this dataset, all models have the same accuracy (83.33%), but it seems that more data is needed to determine the optimal model due to the small data size.

Conclus

Appendix



Thank you!

