

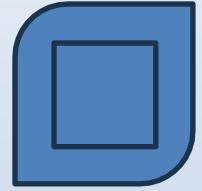
IBM Data Science Capstone Project: Space X

Duyen Nguyen

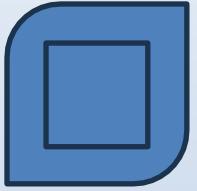
June 26th , 2023



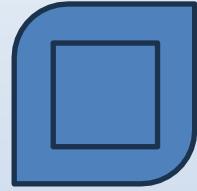
Outline



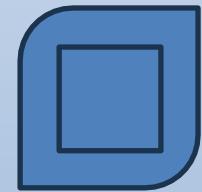
EXECUTIVE
SUMMARY



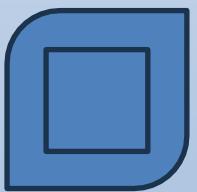
INTRODUCTION



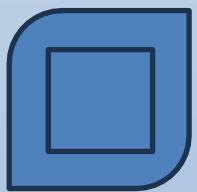
METHODOLOGY



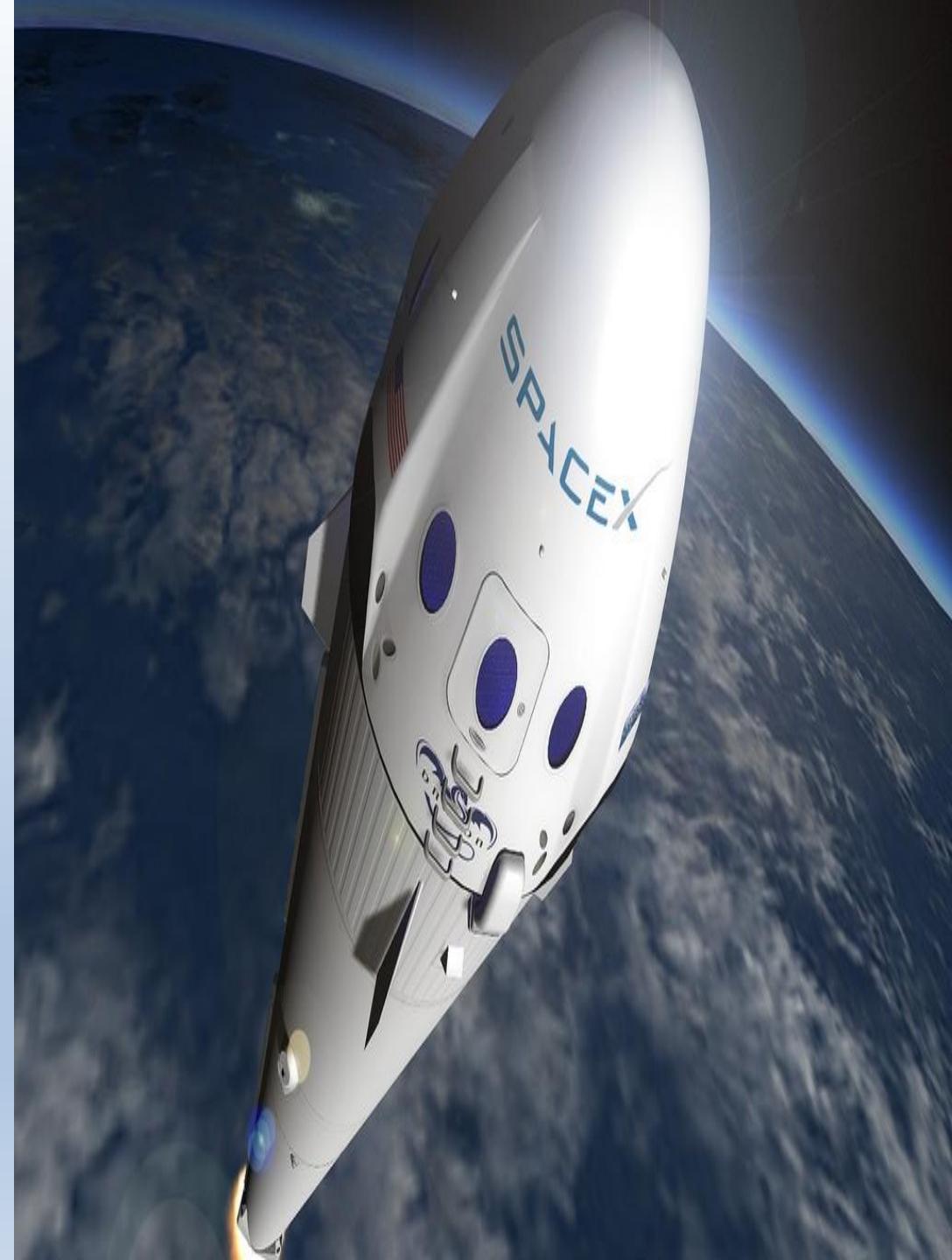
RESULTS



CONCLUSION



APPENDIX



Executive Summary

Summary of methodologies

- Data collection
- Data wrangling
- EDA with data visualization
- EDA with SQL
- Building an Inter EDA results
- Interactive Analytics
- Predictive Analysis results active map with folium
- Building a Dashboard with Plotly Dash
- Predictive Analysis - classification

Summary of all results -

- EDA results
- Interactive Analytics
- Predictive Analytics



Introduction

Project background and context.

- SpaceX advertises Falcon 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

Problems to be addressed

- What influences if the rocket will land successfully?
- What is the effect of each relationship of rockets variables on outcomes?
- What are the Conditions which will aid SpaceX to achieve the best results?

Section 1

Methodology

Methodology

Executive Summary

Data collection methodology:

- SpaceX Rest API
- Web Scrapping from [wikipedia](#)

Perform data wrangling (Transforming data for Machine Learning)

- One Hot Encoding data fields for Machine learning and dropping irrelevant columns

Perform exploratory data analysis (EDA) using visualization and SQL

- Plotting :Scatter and bar graphs to show relationship between variables and to show patterns of data

Perform interactive visual analytics:

- Using Folium and Plotly Dash visualizations

Perform predictive analysis using classification models

- Build, tune, evaluate classification models

Data Collection

The following datasets were collected:

- SpaceX launch data that is gathered from the SpaceX Rest API
- This API will give us data about launches including information about the rocket used, payload delivered, launch specifications, and landing outcome.
- The SpaceX REST API endpoints, or URL, starts with api.spacexdata.com/v4/.
- Another popular data source for obtaining Falcon 9 Launch data is Web Scrapping Wikipedia using BeautifulSoup.



Data Collection - SpaceX API

Data collection-
SpaceX REST API



1. Getting Response from API



2. Converting Response to a .json file



3. Apply custom functions to clean data



4. Assign list to dictionary then dataframe



5. Filter dataframe and export to flat file
(.csv)

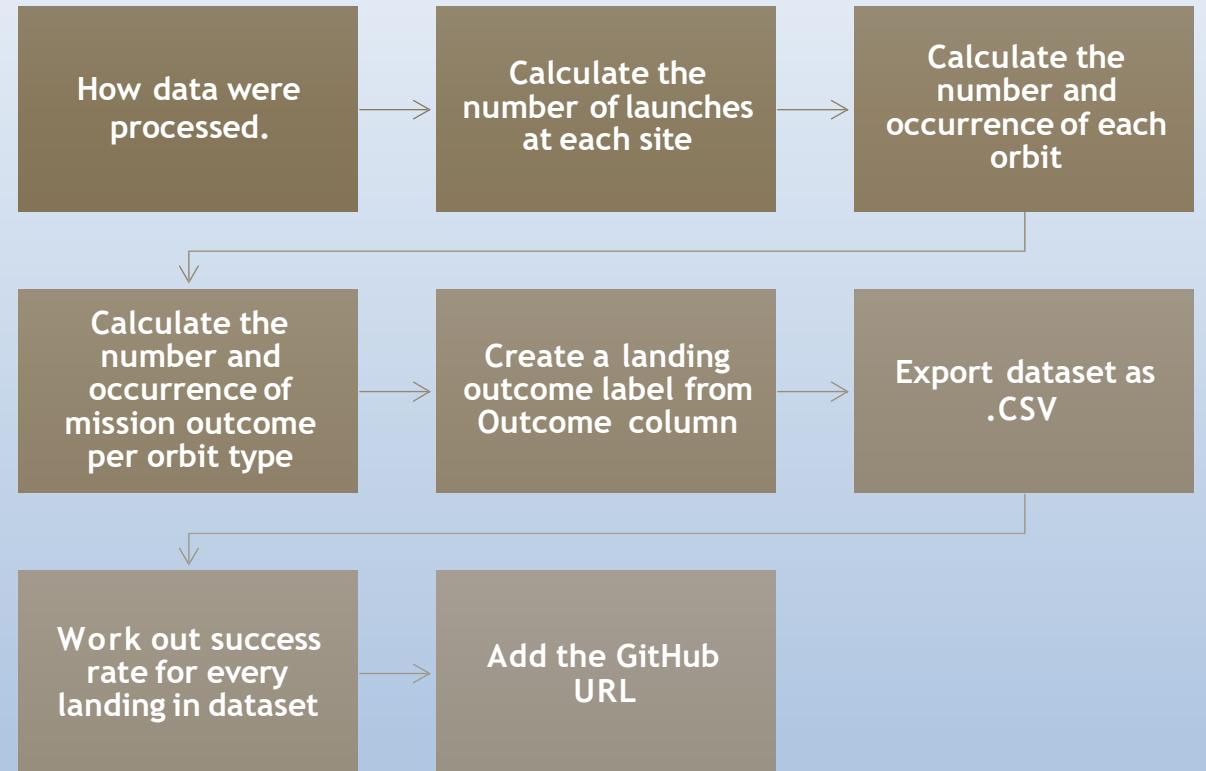
Connect to 

Data Collection - Scrapping



DATA WRANGLING

Connect to  GitHub



EDA with Data Visualization

1. Scatter Graphs: Scatter plots show how much one variable is affected by another. The relationship between two variables is called their correlation. Scatter plots usually consist of a large body of data.

- Flight Number VS. Payload
- Mass Flight Number VS.
- Launch Site Payload VS. Launch Site
- Orbit VS. Flight Number Payload VS. Orbit Type
- Orbit VS. Payload Mass

2. Bar Graph: A bar diagram makes it easy to compare sets of data between different groups at a glance. The graph represents categories on one axis and a discrete value in the other. The goal is to show the relationship between the two axes. Bar charts can also show big changes in data over time.

- Mean VS. Orbit

3. Line Graph: Line graphs are useful in that they show data variables and trends very clearly and can help to make predictions about the results of data not yet recorded

- Success Rate VS. Year

EDA WITH SQL

Summary of the SQL queries you performed

- 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 4000 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 descending order.

Build an Interactive Map with Folium

- Folium makes it easy to visualize manipulated data in python on an interactive leaflet map. We use the latitude and longitude coordinates for each launch site and added a Circle Marker around each launch site with a label of the name of the launch site. Also, it is easy to visualize the number of success and failure for each launch site with green and Red markers on the map

Map Objects	Code	Result
Map Marker	folium.Marker()	Map object to make a mark to map
Icon Marker	folium.icon()	Create an icon on map
Circle Marker	folium.Circle()	Create a circle where marker is being placed
Polyline	folium.Polyline()	Create a line between points.
Marker Cluster Object	MarkerCluster()	This is a good way to simplify a map containing many markers having the same coordinate.
AntPath	Folium.plugins.AntPath()	Create an animated line between points

Build a Dashboard with Plotly Dash

Pie chart shows the total success for all sites / by certain launch site

Scatter graph shows the correlation between payload and success for all sites

Map Objects	Code	Result
Dash and its components	Import dash, import dash_html_components as html	Plotly python's leading data viz and UI libraries. With dash open source, Dash apps run on your local laptop or server. The dash core components library contain a set of higher-level components like slider, graph, dropdown and tables. Dash provides all html tags.
Pandas	Import pandas as pd	Fetching values from CSV and creating a dataframe
Plotly	Import plotly.express as px	Plot the graphs with interactive plotly library
Dropdown	dcc Dropdown()	Create a dropdown for launch sites
Rangeslider	dcc RangeSlider()	Create a rangeslider for payload mass range selection
Pie chart	Px.pie()	Creating the pie graph for success percentage display
Scatter chart	Px.scatter()	Creating the scatter graph for success correlation display

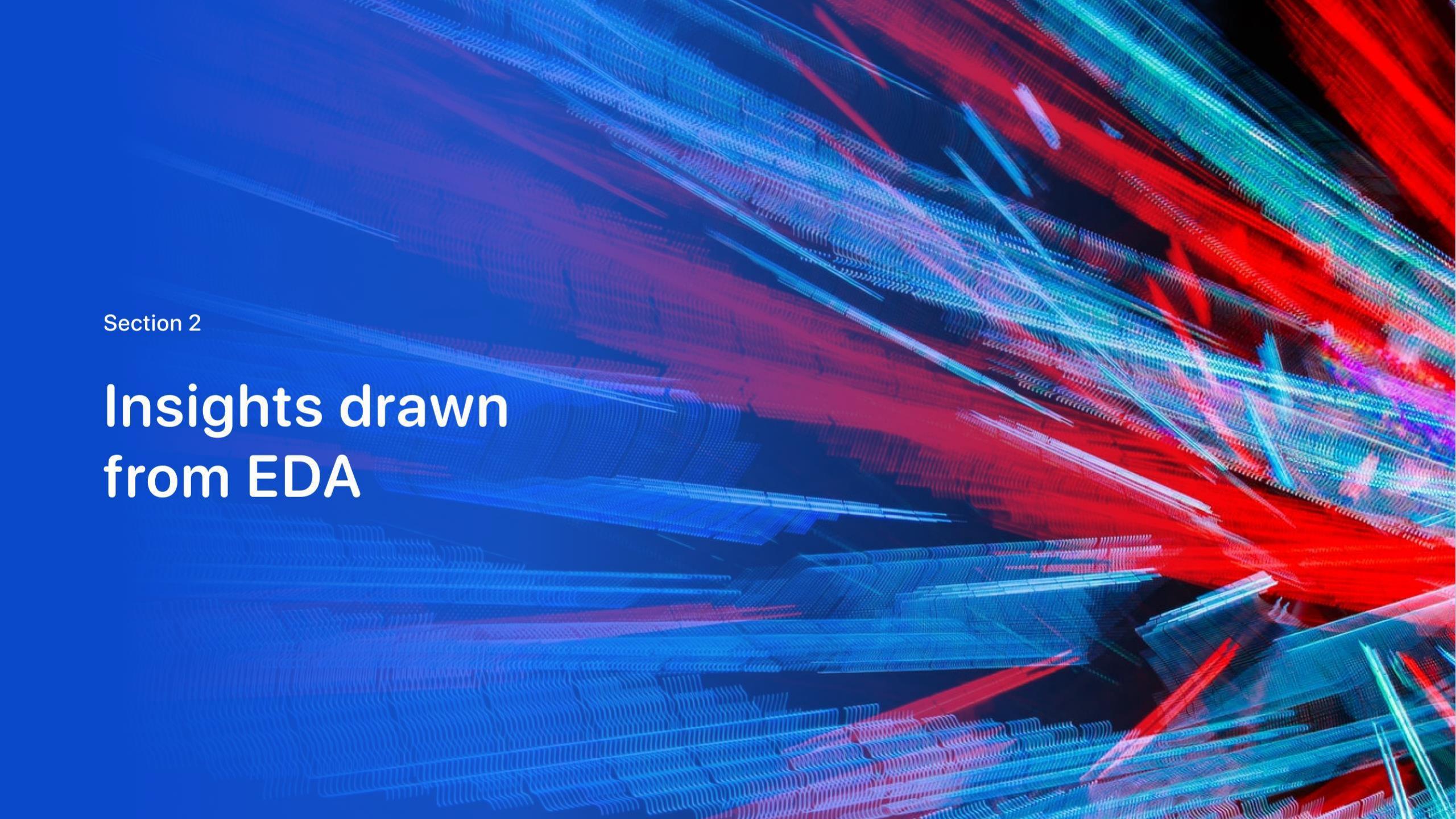
Predictive Analysis (Classification)

- **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 we have
 - Decide which type of machine learning algorithms we want to use
 - Set our parameters and algorithms to GridSearchCv
 - Fit our datasets into the GridSearchCv objects and train our model
- **Evaluating Model**
 - Check accuracy for each model
 - Get tuned hyperparameters for each type of algorithms
 - Plot confusion Matrix
- **Improving Model**
 - Feature Engineering
 - Algorithm Tuning
- **Finding the best performing classification Model**
 - The model with the best accuracy score wins the best performing model





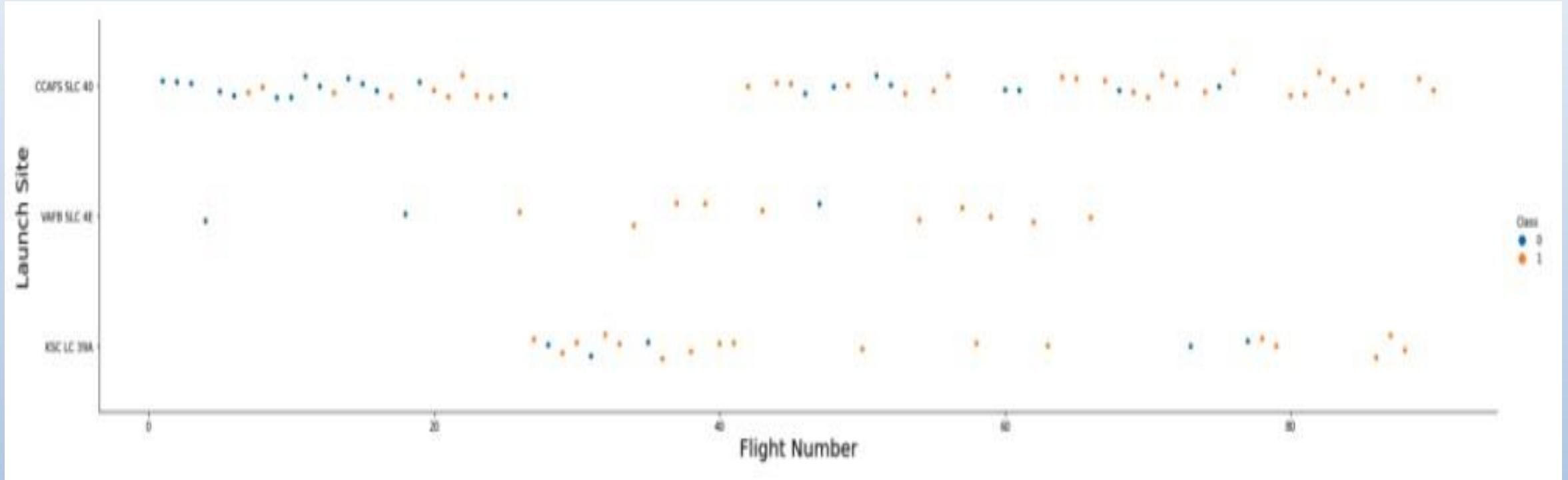
- Exploratory data analysis results
- Interactive analytics demo in screenshots
- Predictive analysis results

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

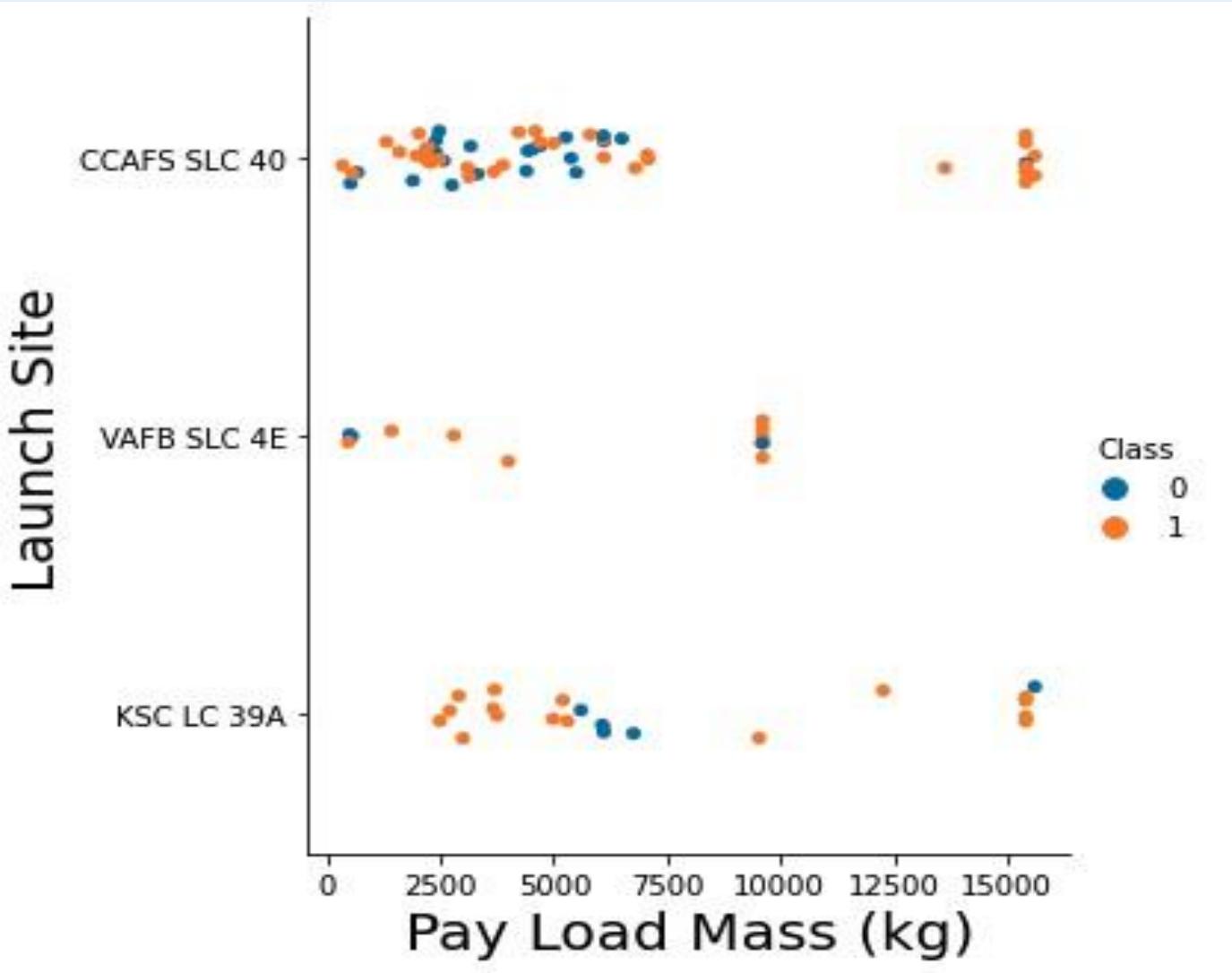
Insights drawn from EDA

Flight Number vs. Launch Site



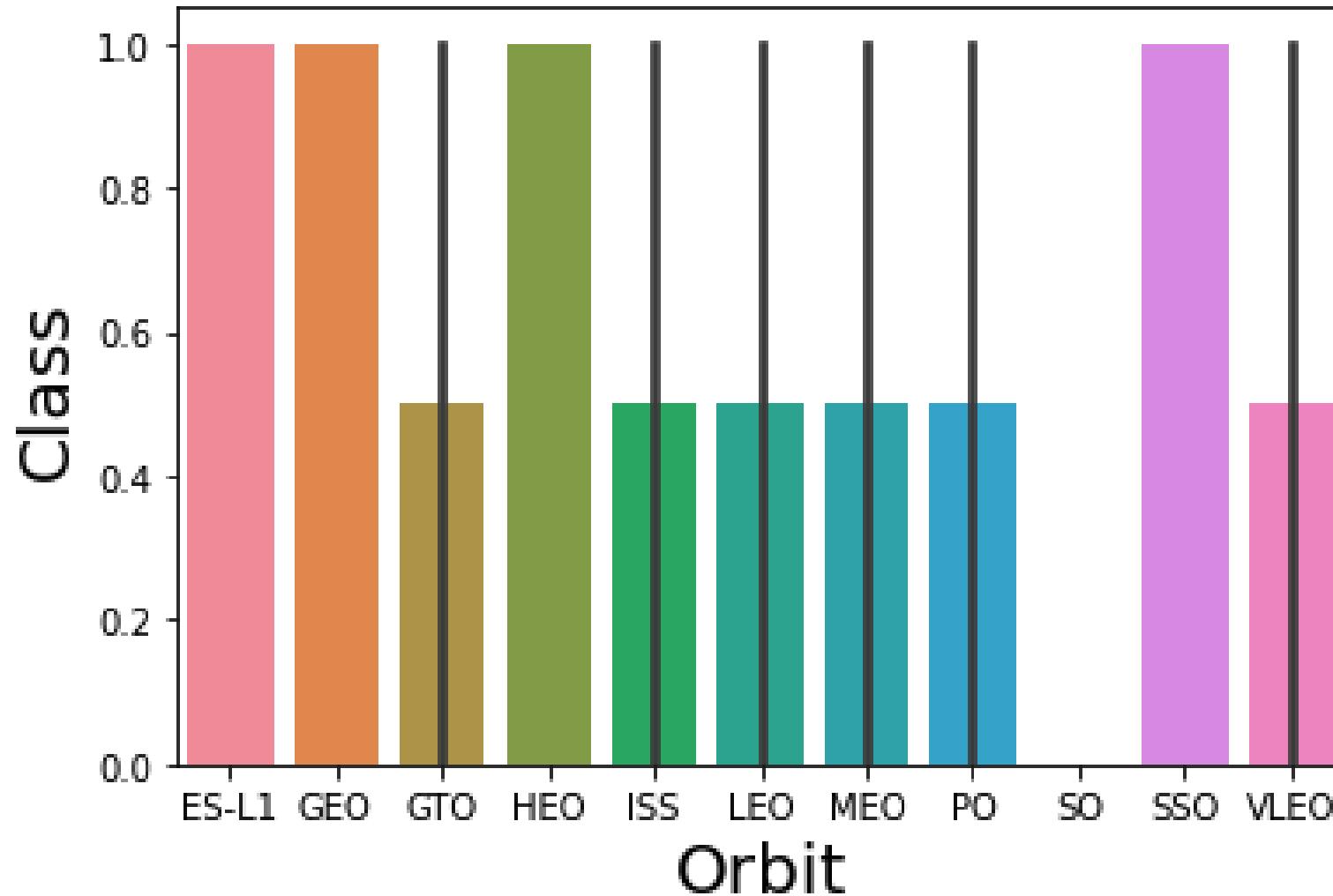
The more flights at a launch site the greater the success rate at a launch site.

Payload vs. Launch Site



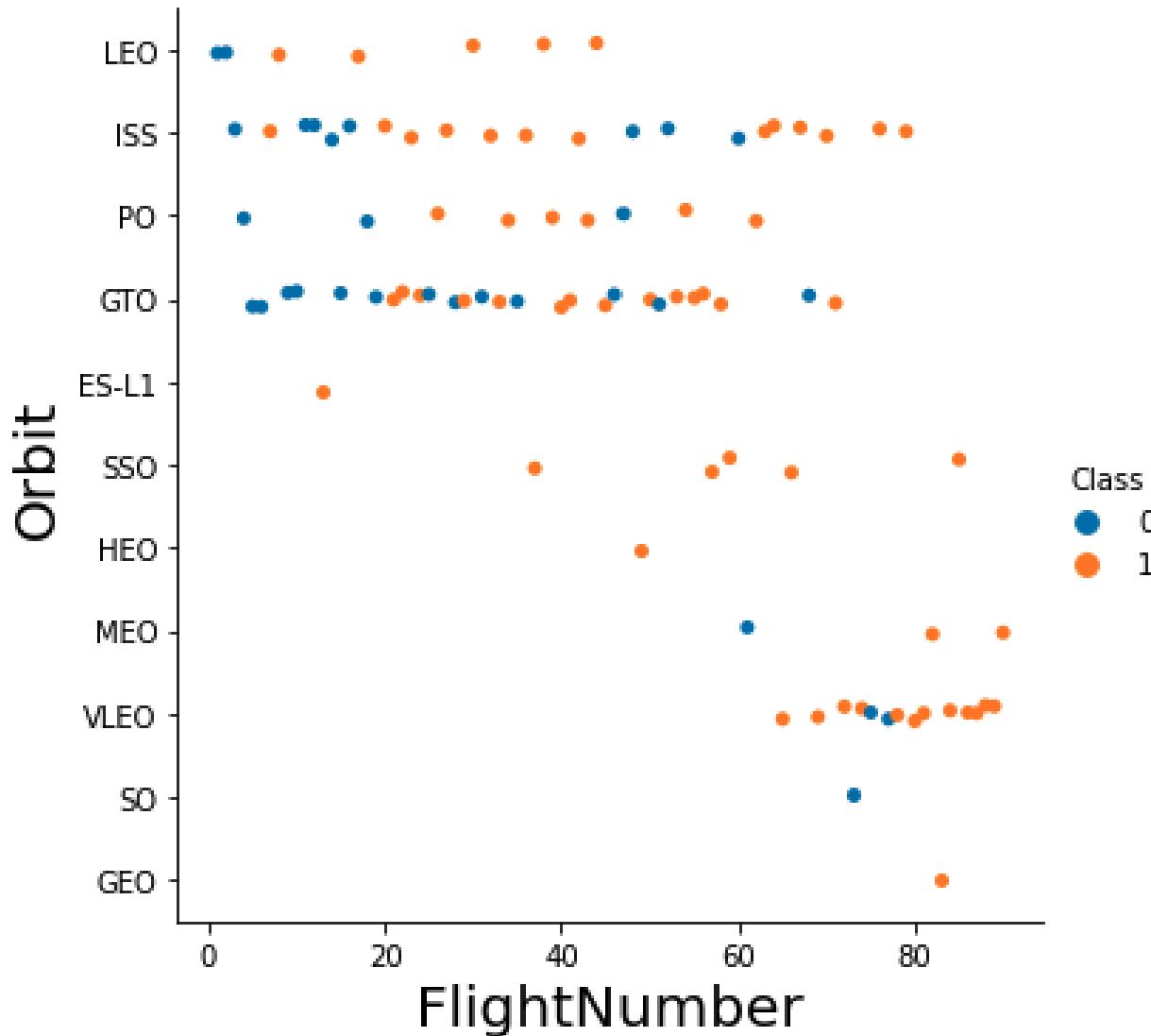
- The more the payload mass for Launch Site CCAFS SLC 40 the higher the success rate for the Rocket. There is not quite a clear pattern to be found using this visualization to make a decision if the Launch Site is dependent on Pay Load Mass for a success launch.

Success Rate vs. Orbit Type



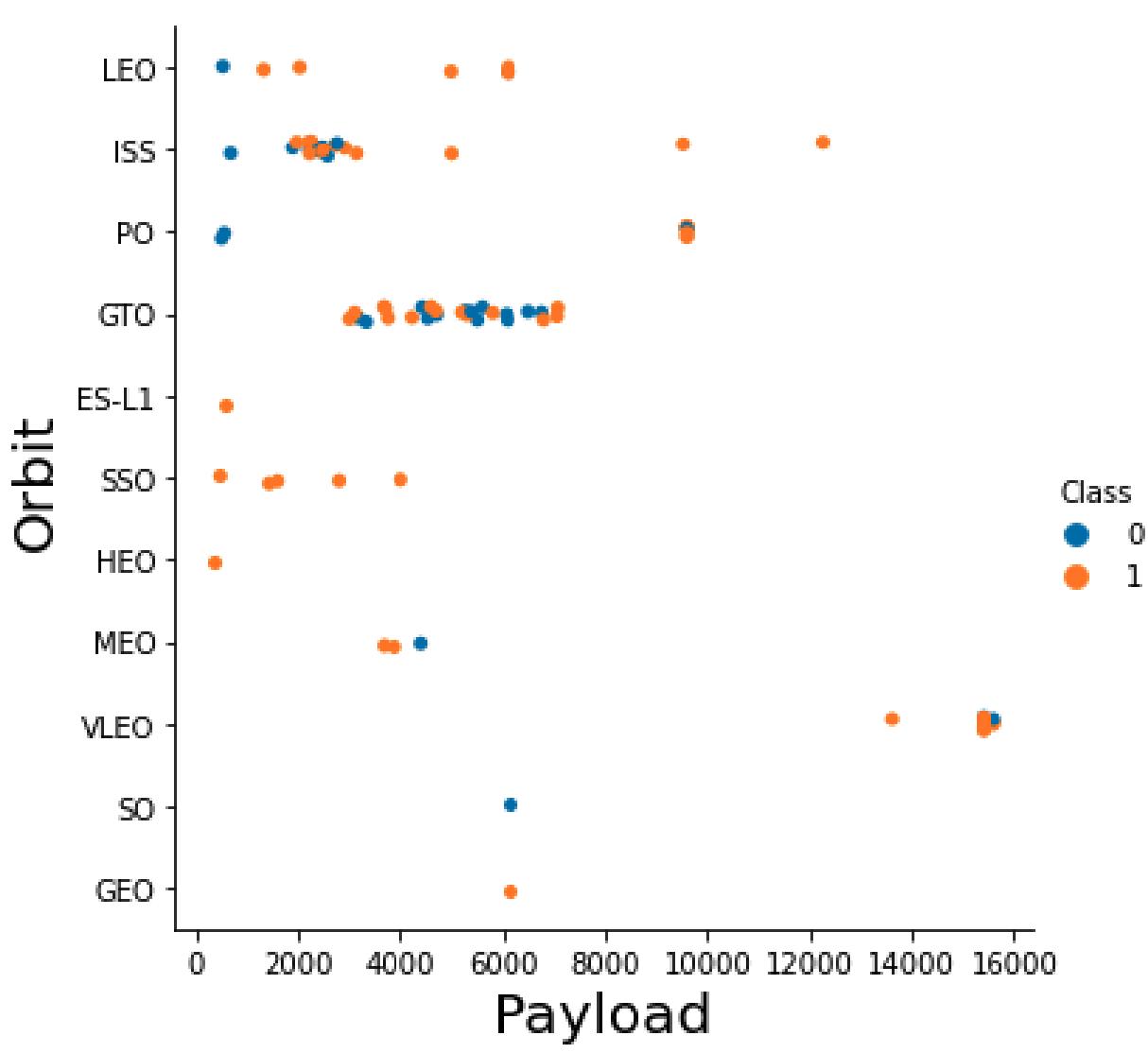
- The Orbit
- GEO, HEO, SSO, ES-L1 has the best Success Rate

Flight Number vs. Orbit Type



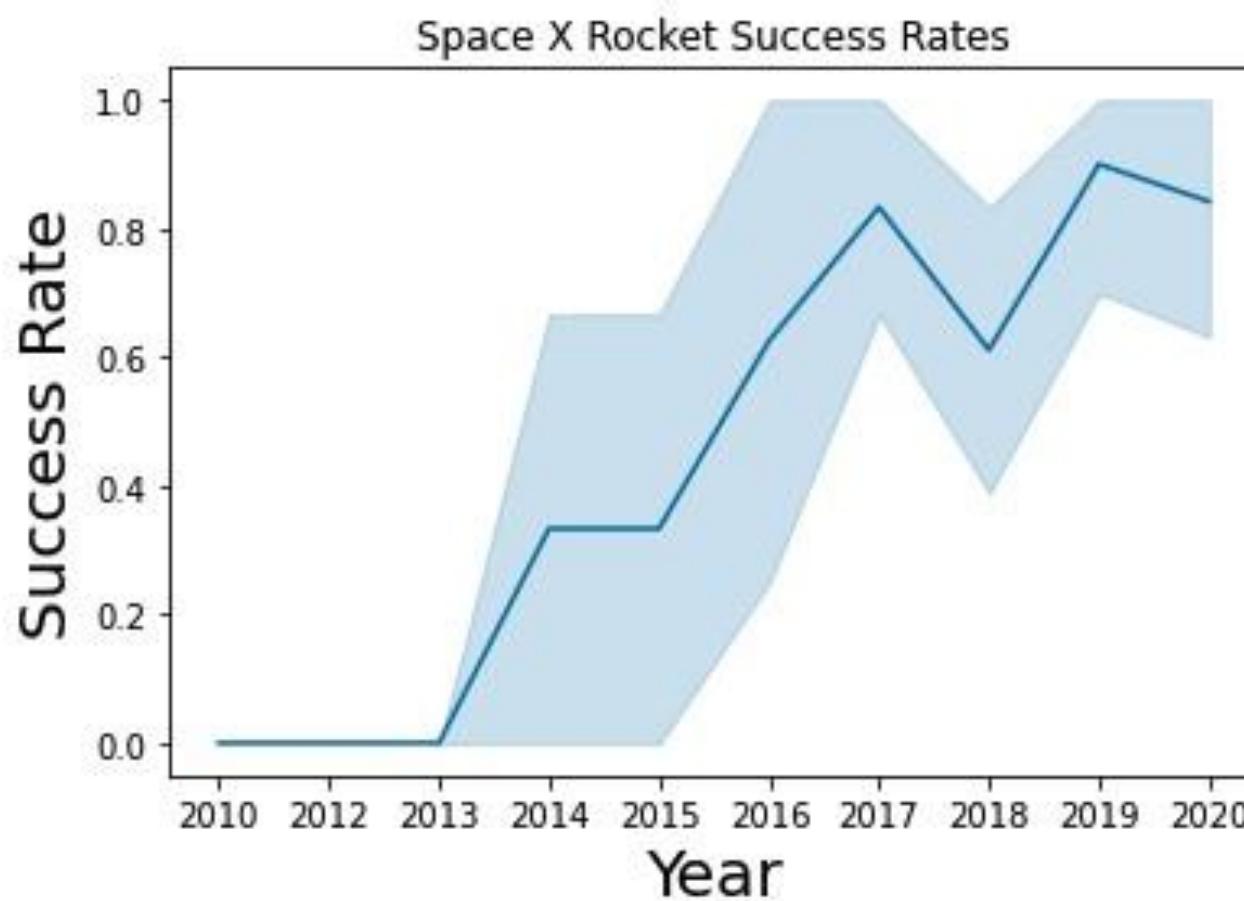
In the LEO orbit the Success appears related to the number of flights, on the other hand, there seems to be no relationship between flight number when in GTO orbit.

Payload vs. Orbit Type



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

Launch Success Yearly Trend



- As we can see that the success rate since 2013 kept increasing till 2020

EDA WITH SQL

All Launch Site Names

- Sql Query

```
*sql select distinct(LAUNCH_Site) from SPACEXTBL
```



Launch_Site

CCAFS LC-40

VAFB SLC-4E

KSC LC-39A

CCAFS SLC-40

Description: Using the word distinct in the query will pull the unique values for the launch Site column from the table SPACEXTBL

Launch Site Names Begin with 'CCA'

- Sql Query

```
%sql select * from SPACEXTBL where LAUNCH_SITE like 'CCA%' limit 5
```

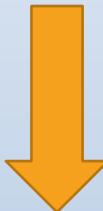
- Description: Using keyword “Limit 5” in the query will fetch 5 records from table spacex, condition LIKE keyword with wild card “CCA%”. The percentage in the end suggest that the launch_site name must start with CCA.

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

- SQL Query

```
%sql select sum(PAYLOAD_MASS__KG_) from SPACEXTBL where CUSTOMER = 'NASA (CRS)';
```



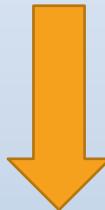
sum(PAYLOAD_MASS_KG_)
45596

Description: Using function SUM summate the total in the column PAYLOAD_MASS_KG and the WHERE clause filters the dataset to only perform calculations on customer NASA (CRS)

Average Payload Mass by F9 v1.1

- **Sql query**

```
%sql select avg(PAYLOAD_MASS_KG_) from SPACEXTBL where BOOSTER_VERSION LIKE '%F9 v1.1';
```



avg(PAYLOAD_MASS_KG_)
2928.4

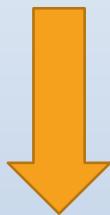
Description:

Using the function AVG works out the average in the column PAYLOAD_MASS_KG

The WHERE clause filters the dataset to only perform calculations on Booster_version f9 v1.1

First Successful Ground Landing Date

- SQL Query:
- select MIN(Date) SLO from tblSpaceX where Landing_Outcome = "Success (drone ship)"



Date which first Successful landing outcome in drone ship was acheived.

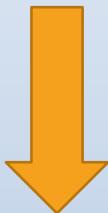
0

06-05-2016

- Description: The function MIN works out the minimum date in the column date while the WHERE clause filters the dataset to only perform calculations on Landing_outcome success

Successful Drone Ship Landing with Payload between 4000 and 6000

- SQL Query:
- `select Booster_Version from tblSpaceX where Landing_Outcome = 'Success (ground pad)' AND Payload_MASS_KG_ > 4000 AND Payload_MASS_KG_ < 6000`



Date which first Successful landing outcome in drone ship was acheived.

0	F9 FT B1032.1
1	F9 B4 B1040.1
2	F9 B4 B1043.1

Description: Selecting only Booster_Version. WHERE clause filters AND clause specifies additional filter.

[Github url to notebook](#)

Total Number of Successful and Failure Mission Outcomes

- SQL Query:
- %sql select count(MISSION_OUTCOME) from SPACEXTBL where MISSION_OUTCOME = 'Success' or MISSION_OUTCOME = 'Failure (in flight)'



count(MISSION_OUTCOME)
99

Boosters Carried Maximum Payload

- Sql Query:

```
sql select BOOSTER_VERSION from SPACEXTBL where PAYLOAD_MASS_KG_ = (select max(PAYLOAD_MASS_KG_) from SPACEXTBL)
```



Booster_Version
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

Description: Using the function MAX works out the maximum payload in the column PAYLOAD_MASS_KG_ in the sub query and WHERE clause filters Booster Version which had that maximum payload.

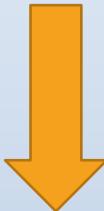
2015 Launch Records

- SQL Query
- %sql SELECT BOOSTER_VERSION, LAUNCH_SITE FROM SPACEXTBL WHERE DATE LIKE '2015-%' AND \
- LANDING_OUTCOME = 'Failure (drone ship)';

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

- SQL query:

```
%sql select Count(LANDING_OUTCOME) AS "Rank success count between 2010-06-04 and 2017-03-20" from SPACEXTBL \
where LANDING_OUTCOME like 'Success%' and (DATE between '2010-06-04' and '2017-03-20') order by date desc
```



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

8

Description:

COUNT counts records in column
LANDING_OUTCOME. WHERE filters data with
'%Success%'

The background of the slide is a photograph taken from space at night. It shows the curvature of the Earth's horizon against a dark blue sky. 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

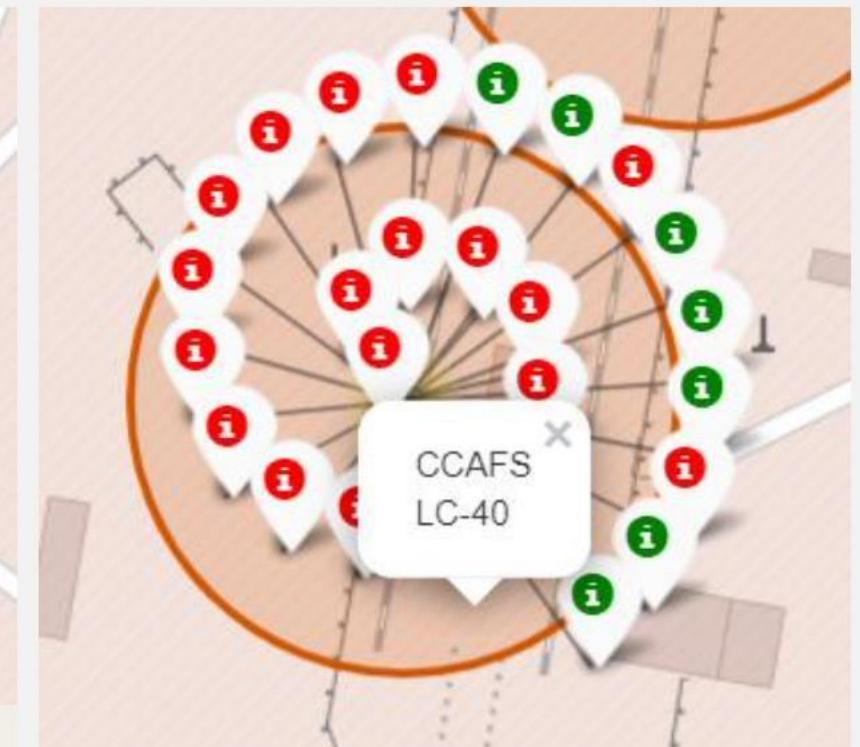
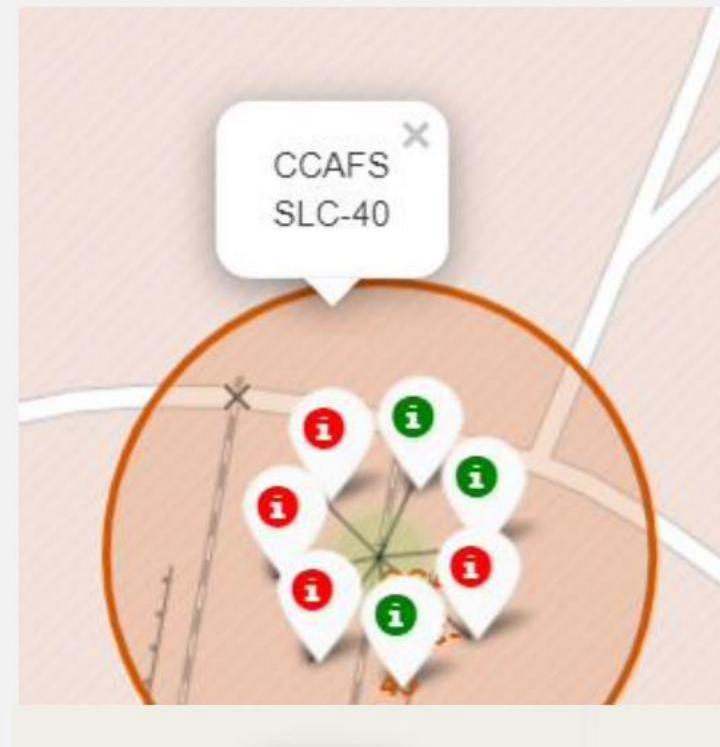
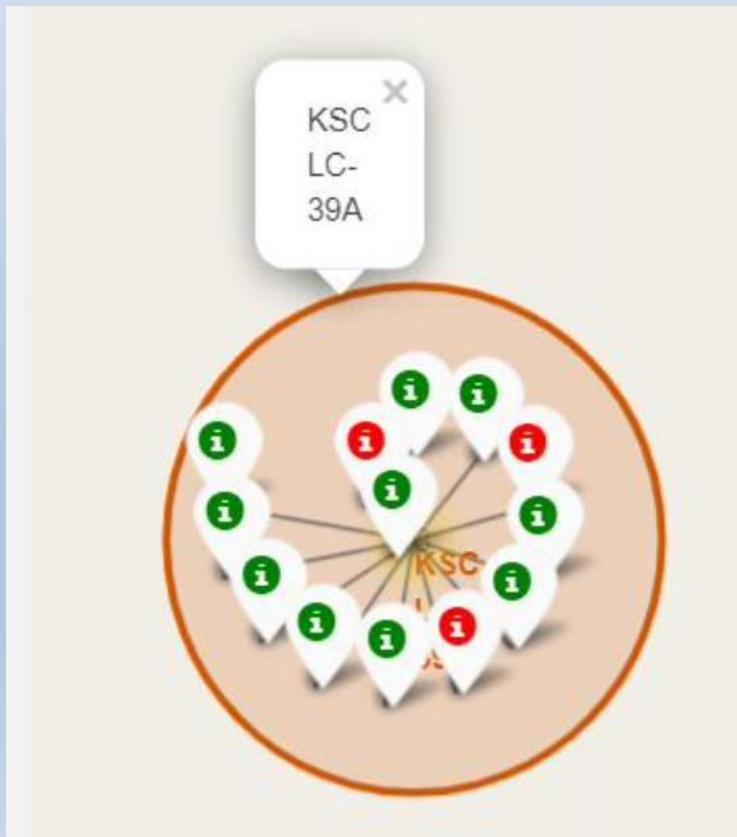
All launch sites global map markers



According to the map:

The SpaceX launch sites are in the united states of America coast.
Florida and California

Colour Labelled Markers

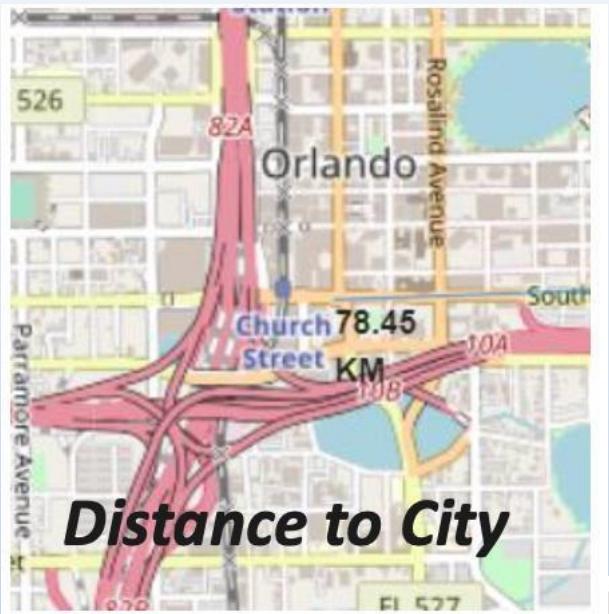


Green Marker shows successful launches and Red Marker shows Failures

Working out launch sites distance to landmark to find trends with haversine formula using CCAFS-SLC-40 as a reference



Distance to coast



Distance to City

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?

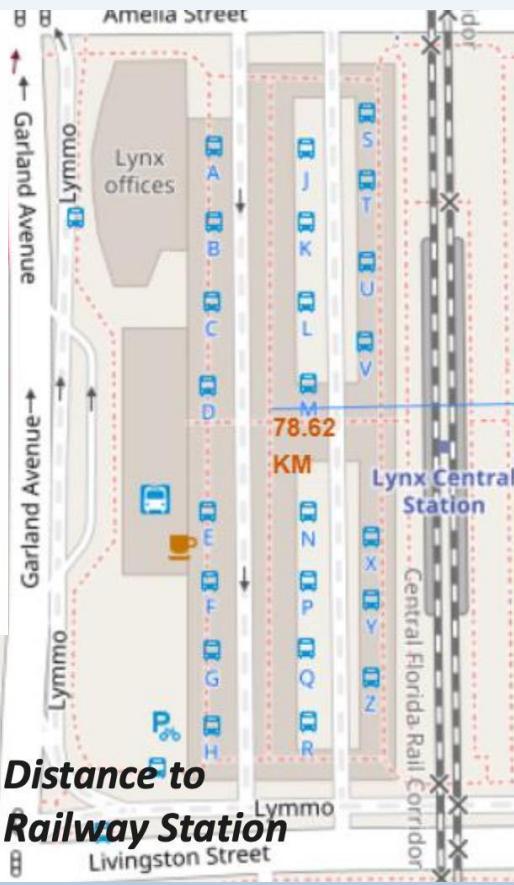
YEs



Distance to Coastline



Distance to closest Highway

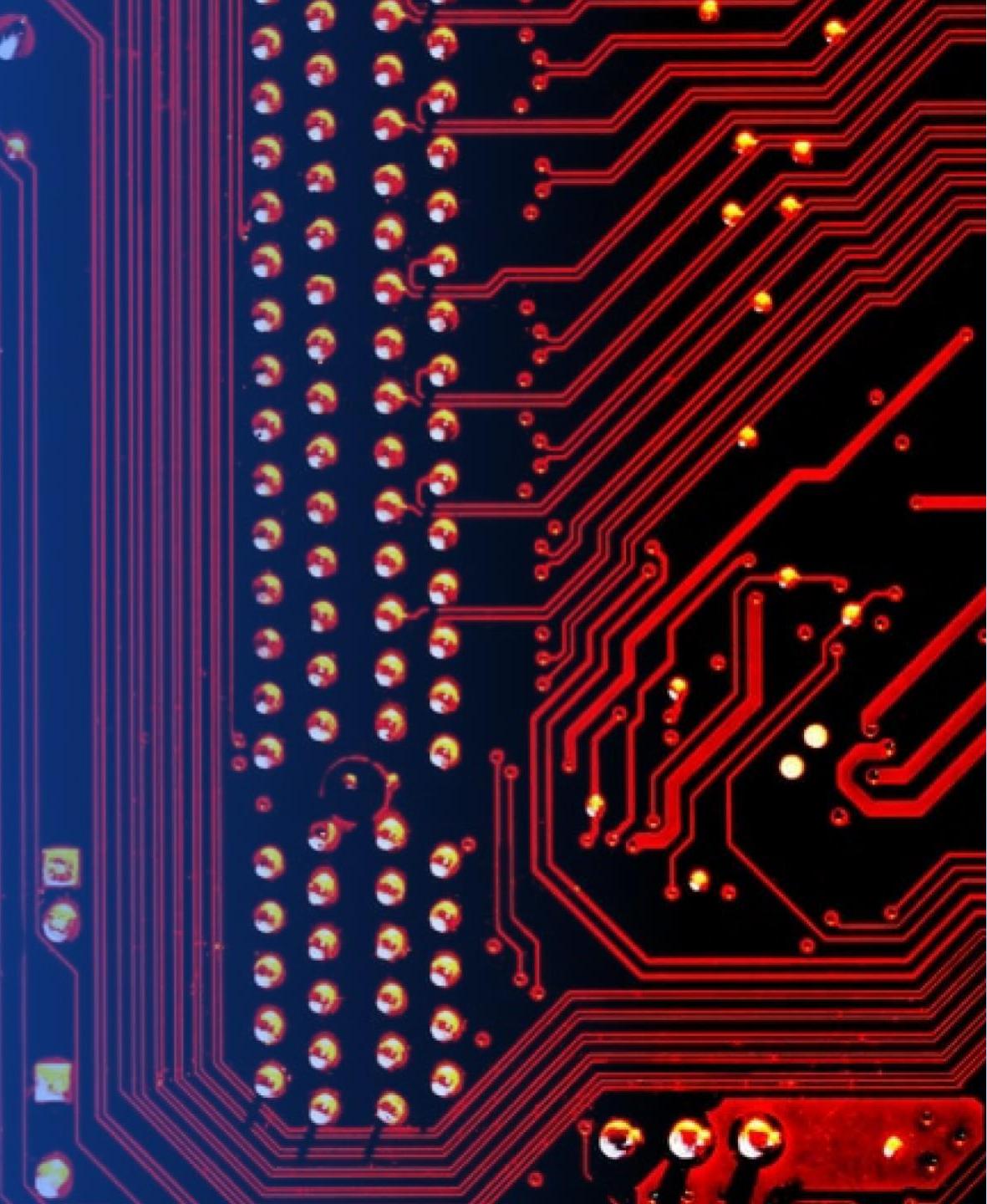


Distance to Railway Station

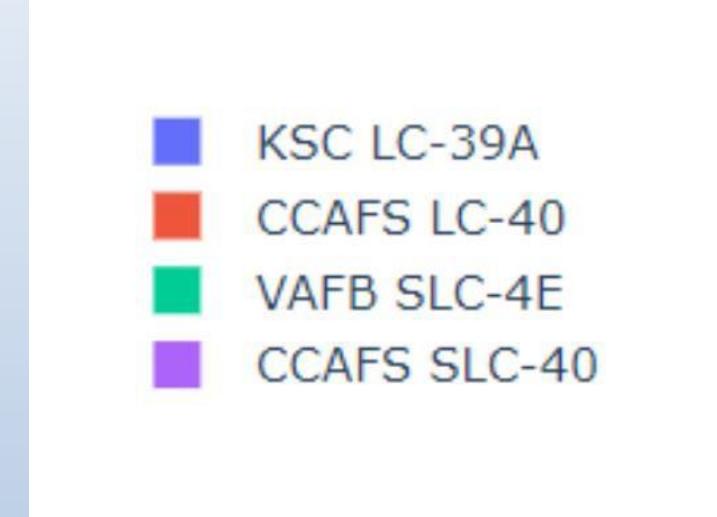
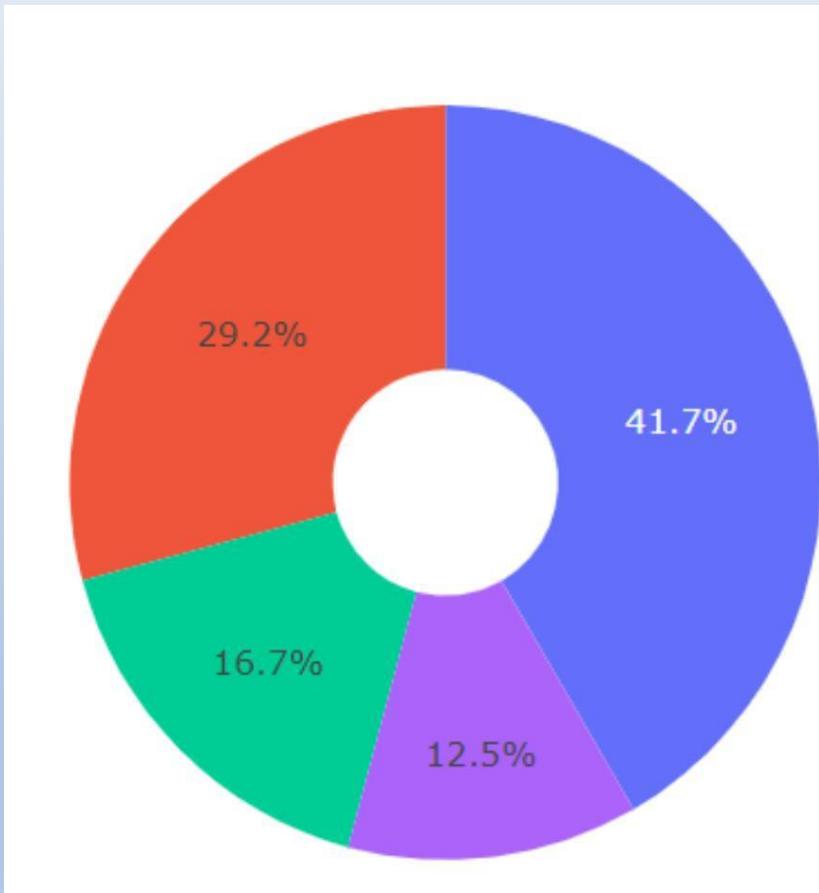
[Github url to notebook](#)

Section 5

Build a Dashboard with Plotly Dash

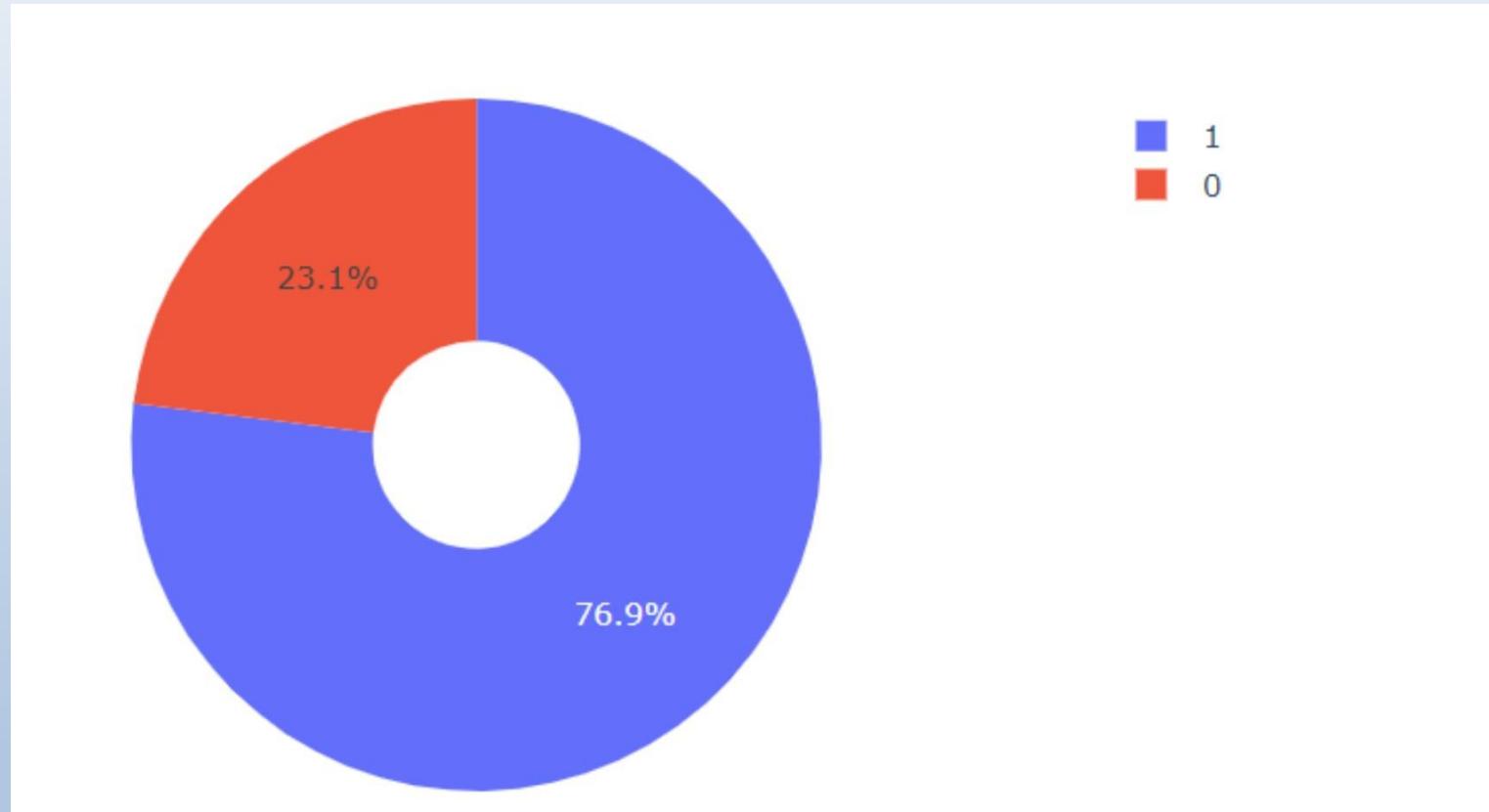


Total success launches by all sites



As we can see from the pie chart
KSC LC-39A had the most successful launches from all the sites

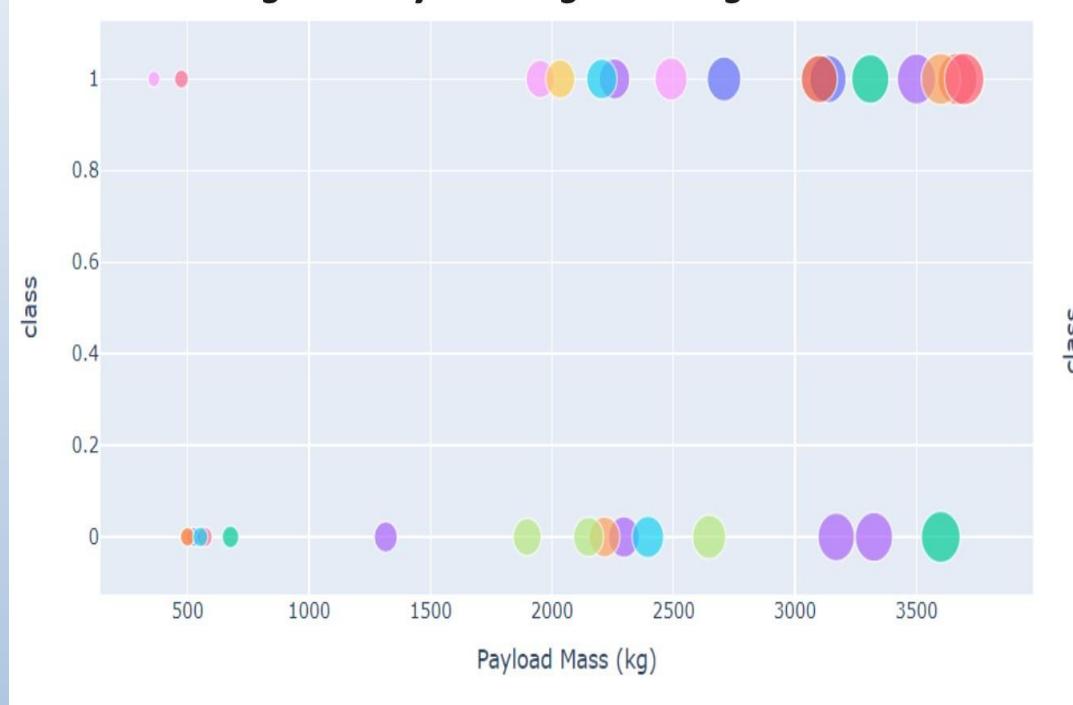
DASHBOARD - Pie chart for the launch site with highest launch success ratio



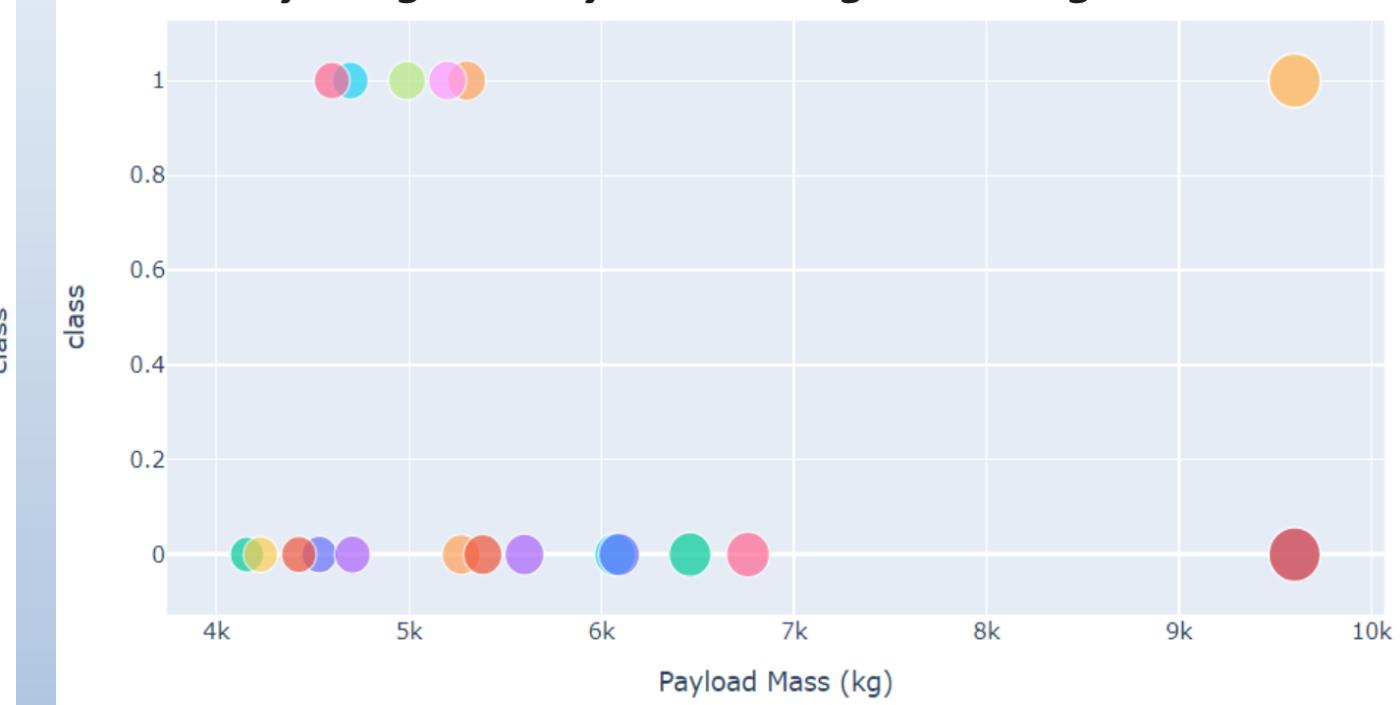
KSC LC-39A achieved a 76.9% success rate while getting a 23.1% failure rate

Dashboard- Payload vs launch outcome scatter plot for all sites, With different payload selected in the range slider

Low weighted payload 0kg - 4000kg



Heavy weighted payload 4000kg - 10000kg



As we can see the success rates for low weighted payloads is higher than the heavy weighted payloads

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

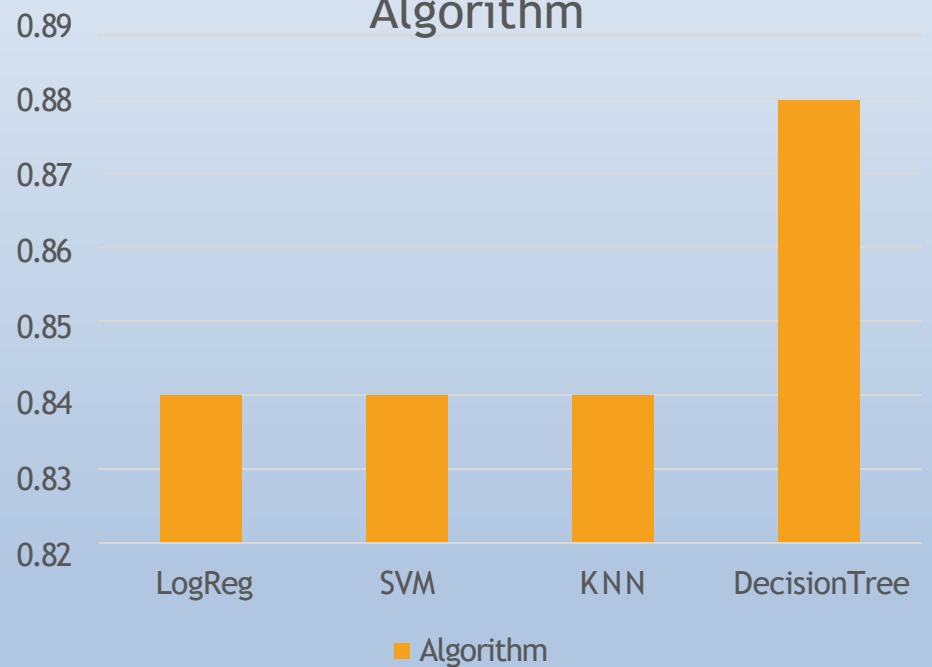
Predictive Analysis (Classification)

Classification Accuracy

The best perform accuracy is Decision tree with a score of 0.88. We trained four model and none of them had anything less than 0.83

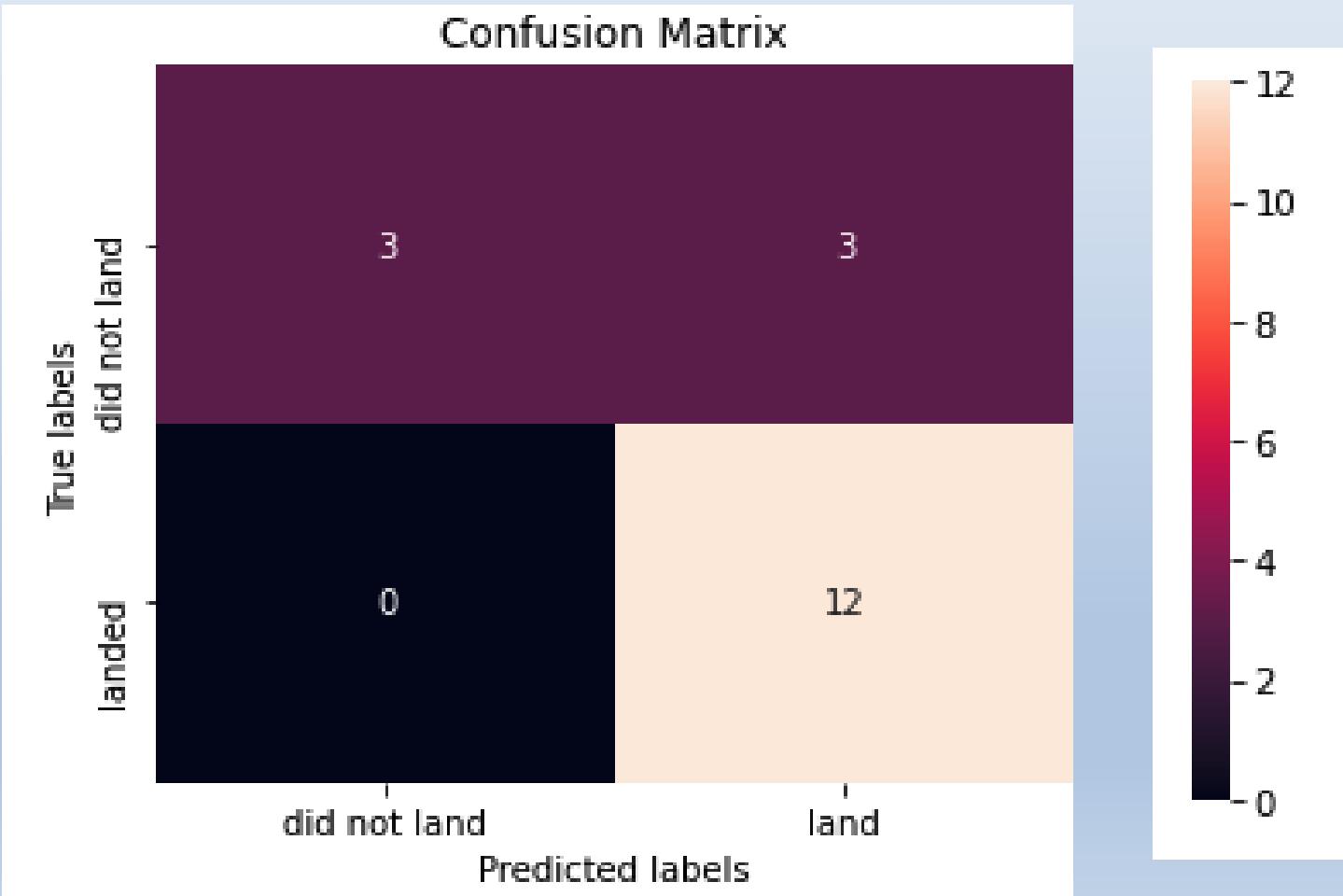
Algorithm	Accuracy	Accuracy on test data
Logistic Regression	0.8464	0.83334
SVM	0.84821	0.83334
KNN	0.8482	0.83334
Decision Tree	0.8892	0.7222

Bar chart showing Accuracy vs Algorithm



Confusion Matrix

All models used have the same confusion matrix.



CONCLUSIONS

Orbits ES-L1, GEO, HEO, SSO has highest success rates

The Success rates for SpaceX launches has been increasing relatively with time, they will eventually perfect the launches.

KSC LC_39A had the most successful launches from all the sites

Low weighted payloads perform better than the heavier payloads

The tree classifier Algorithm is the best machine learning model for this dataset

APPENDIX

Interactive plotly

Python Anywhere

Folium MeasureControl Plugin Tool

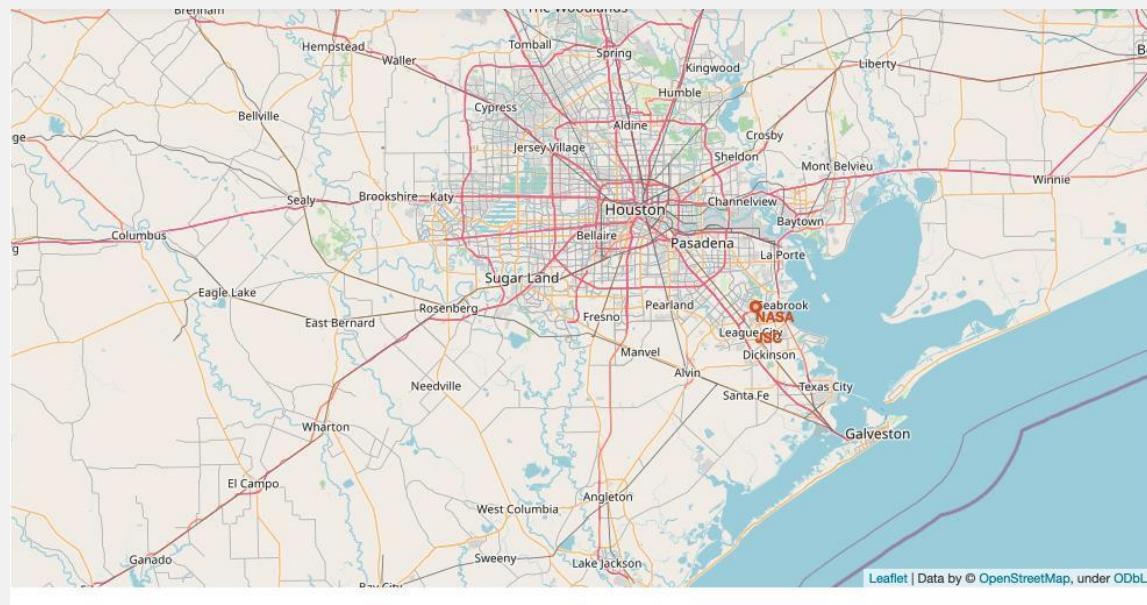
Basic Decision Tree Construction

IBM cognos Visualization Tool

Folium MeasureControl Plugin Tool

I used Folium custom to create custom layer to understand the location of the launch site

```
lines=folium.PolyLine(locations=coordinates, weight=1)
site_map.add_child(lines)
distance = calculate_distance(coordinates[0][0], coordinates[0][1], coordinates[1][0], coordinates[1][1])
distance_circle = folium.Marker(
    [28.56342, -80.56794],
    icon=DivIcon(
        icon_size=(20,20),
        icon_anchor=(0,0),
        html='<div style="font-size: 12; color:#d35400;"><b>%s</b></div>' % "{:10.2f} KM".format(distance),
    )
)
site_map.add_child(distance_circle)
site_map
```



Thank
You

