



IBM Developer  
SKILLS NETWORK

# Winning Space Race with Data Science

<Name>  
<Date>



# 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
  - Build, Train, and Evaluate 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 notebook is  
<https://github.com/yenchun19971028/Space-X/blob/main/jupyter-labs-spacex-data-collection-api.ipynb>

```
In [ ]: spacex_url="https://api.spacexdata.com/v4/launches/past"
In [ ]: response = requests.get(spacex_url)
Check the content of the response
In [ ]: print(response.content)
```

You should see the response contains massive information about SpaceX launches. Next, let's try to discover some more relevant information for this project.

**Task 1: Request and parse the SpaceX launch data using the GET request**

To make the requested JSON results more consistent, we will use the following static response object for this project:

```
In [ ]: static_json_url='https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/
```

We should see that the request was successfull with the 200 status response code

```
In [ ]: response.status_code
```

Now we decode the response content as a Json using `.json()` and turn it into a Pandas dataframe using `.json_normalize()`

```
In [ ]: # Use json_normalize meethod to convert the json result into a dataframe
```

# 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/yenchun19971028/Space-X/blob/main/jupyter-labs-webscraping.ipynb>

```
In [4]: static_url = "https://en.wikipedia.org/w/index.php?title=List_of_Falcon_9_and_Fa
```

Next, request the HTML page from the above URL and get a `response` object

## TASK 1: Request the Falcon9 Launch Wiki page from its URL

First, let's perform an HTTP GET method to request the Falcon9 Launch HTML page, as an HTTP response.

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

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

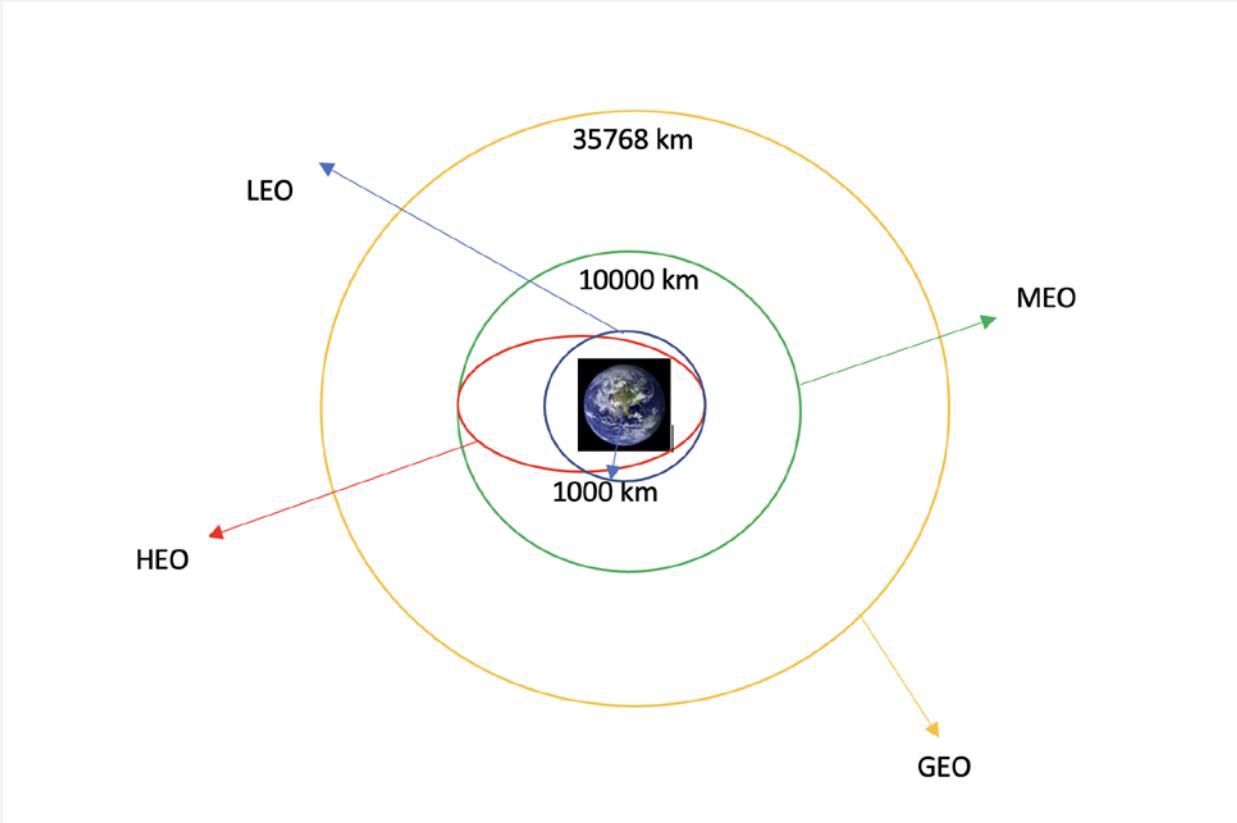
```
In [8]: # Use BeautifulSoup() to create a BeautifulSoup object from a response text content  
bs = BeautifulSoup(response.content,"html.parser")
```

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

```
In [9]: # Use soup.title attribute  
bs.title
```

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

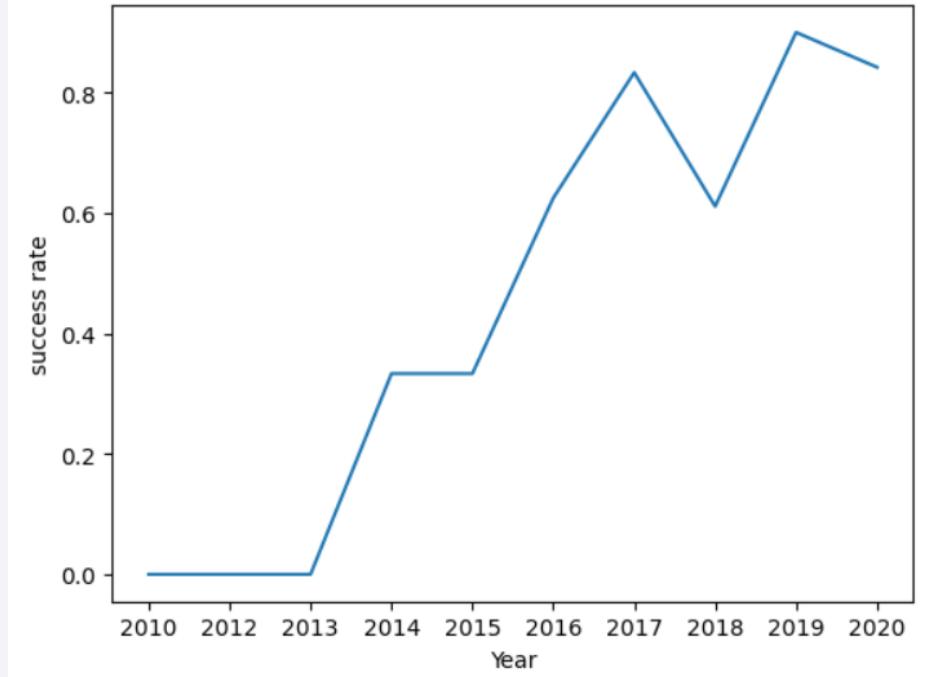
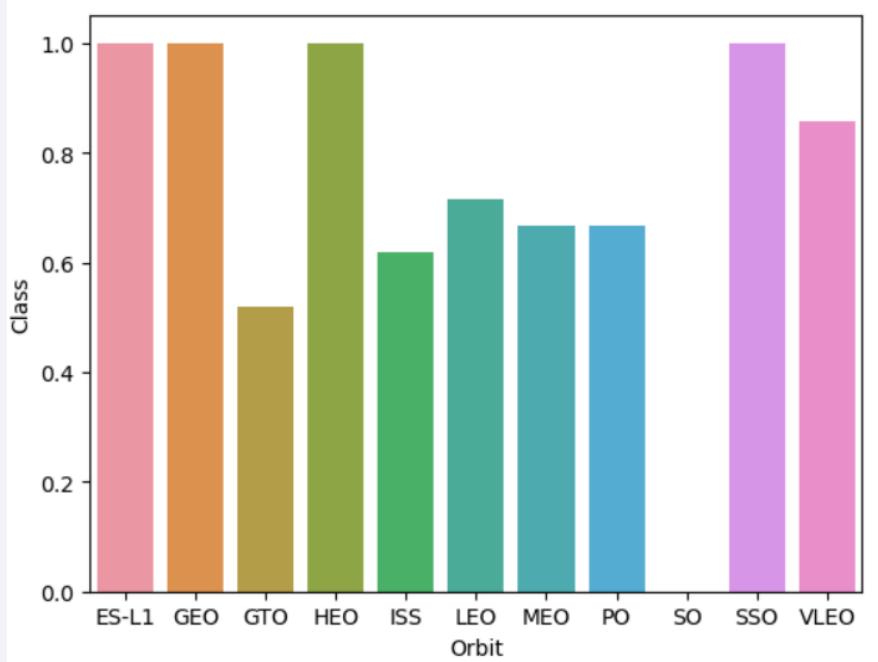
# Data Wrangling



- 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/yenchun19971028/Space-X/blob/main/labs-jupyter-spacex-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/yenchun19971028/  
Space-X/blob/main/jupyter-labs-eda-  
dataviz.ipynb](https://github.com/yenchun19971028/Space-X/blob/main/jupyter-labs-eda-dataviz.ipynb)

# EDA with SQL

---

- We loaded the SpaceX dataset into a PostgreSQL database without leaving the jupyter notebook.
- 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/chuksoo/IBM-Data-Science-Capstone-SpaceX/blob/main/EDA%20with%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.

# Build a Dashboard with Plotly Dash

---

- 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.
- The link to the notebook is [https://github.com/yenchun19971028/Space-X/blob/main/spacex\\_dash\\_app.py](https://github.com/yenchun19971028/Space-X/blob/main/spacex_dash_app.py)

# Predictive Analysis (Classification)

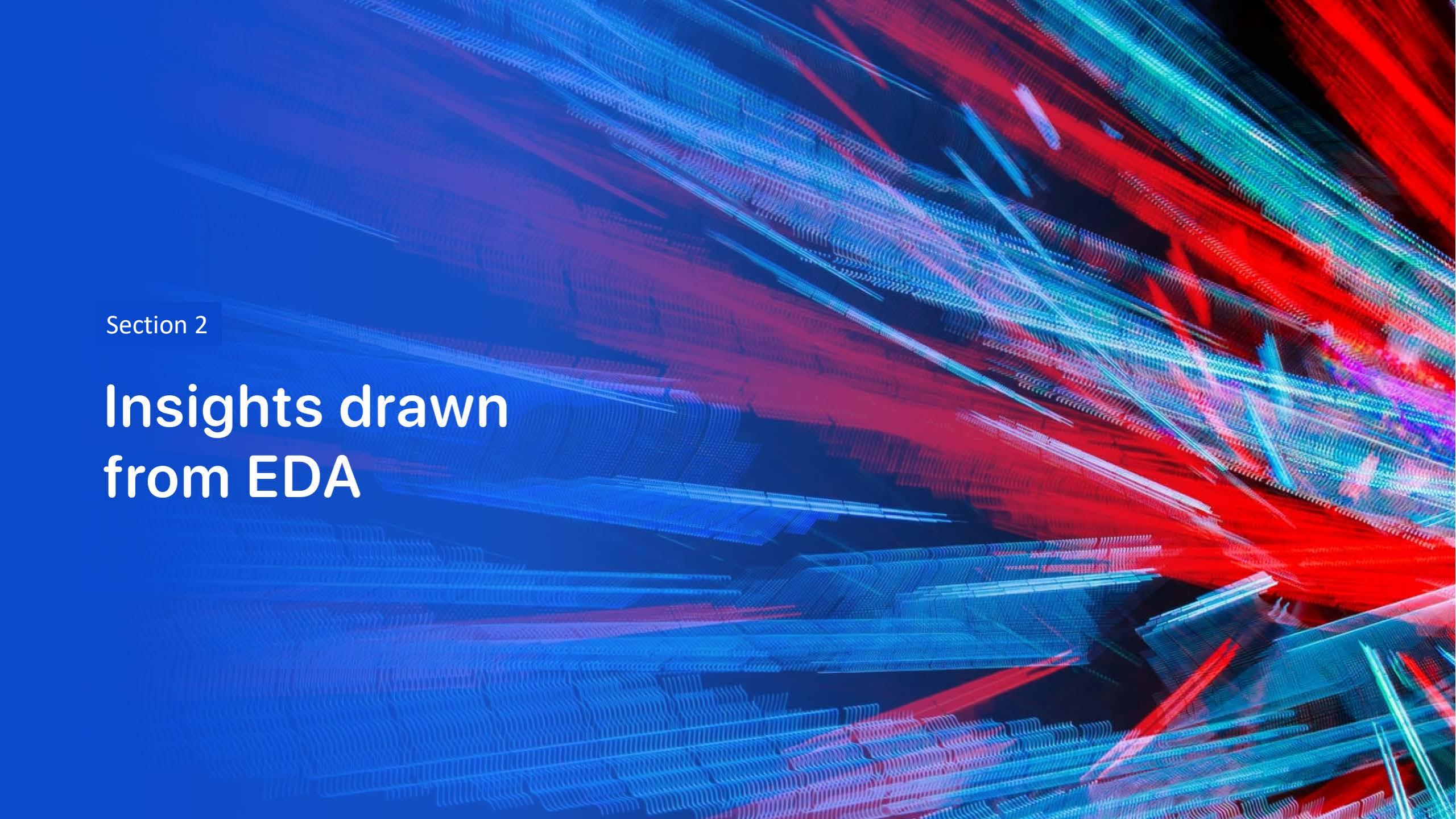
---

- 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/yenchun19971028/SpaceX/blob/main/SpaceX\\_Machine%20Learning%20Prediction\\_Part\\_5.ipynb](https://github.com/yenchun19971028/SpaceX/blob/main/SpaceX_Machine%20Learning%20Prediction_Part_5.ipynb)

# Results

---

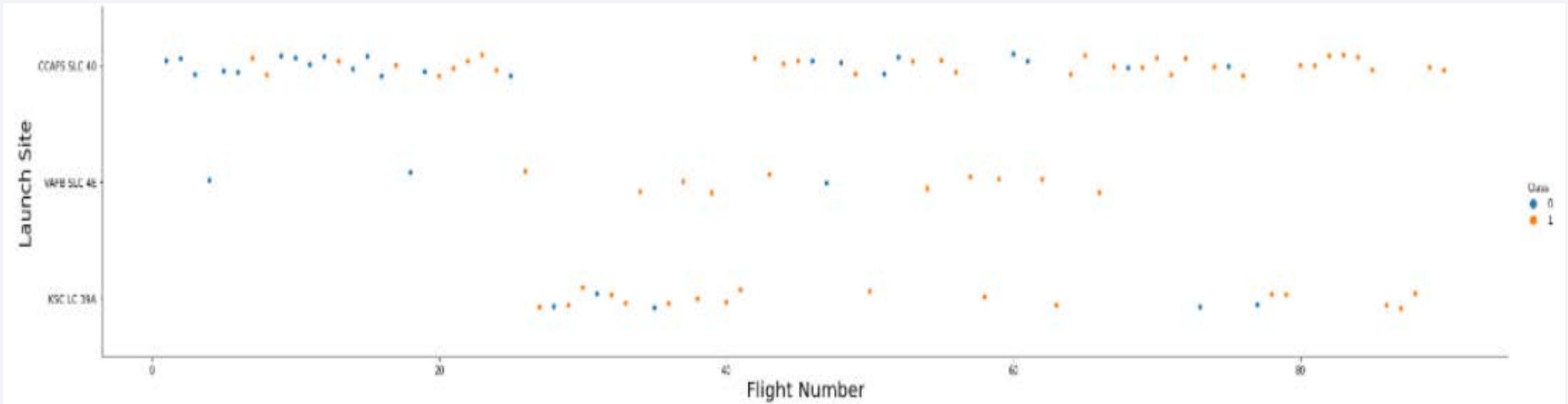
- Exploratory data analysis results
- Interactive analytics demo in screenshots
- Predictive analysis results

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 is more dense and vibrant towards the right side of the frame, while appearing more sparse and blue-tinted on the left. The overall effect is reminiscent of a high-energy particle simulation or a futuristic circuit board.

Section 2

## Insights drawn from EDA

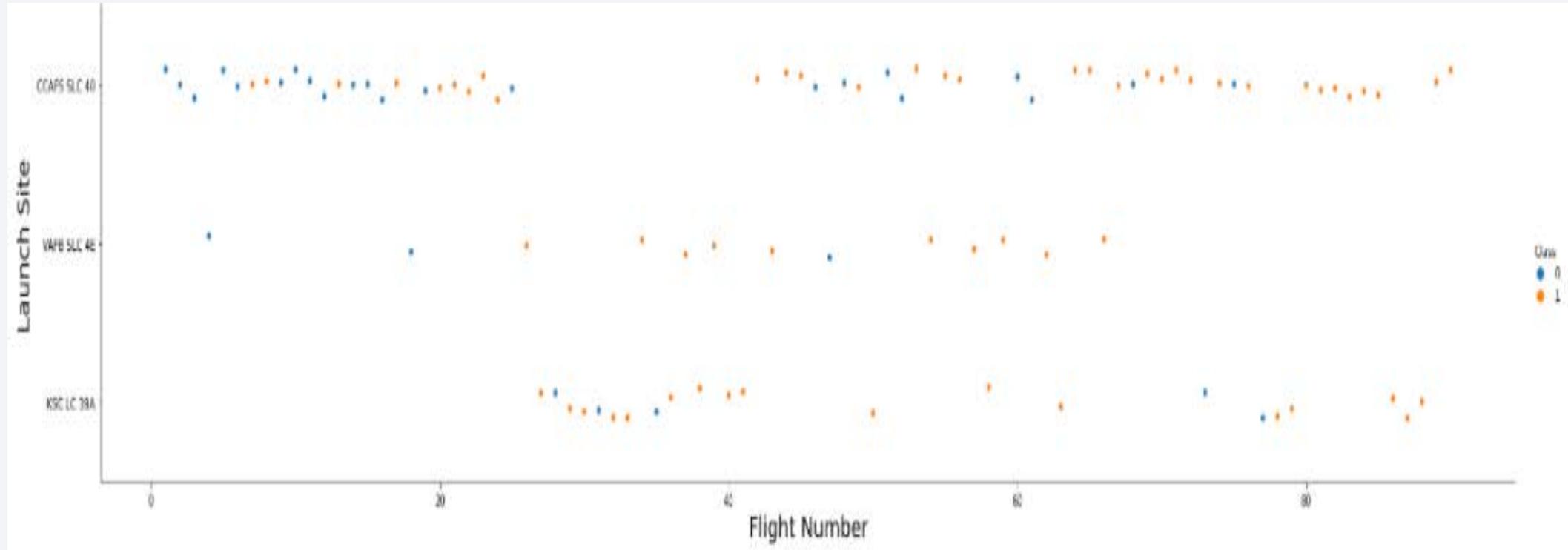
# Flight Number vs. Launch Site



- 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

---

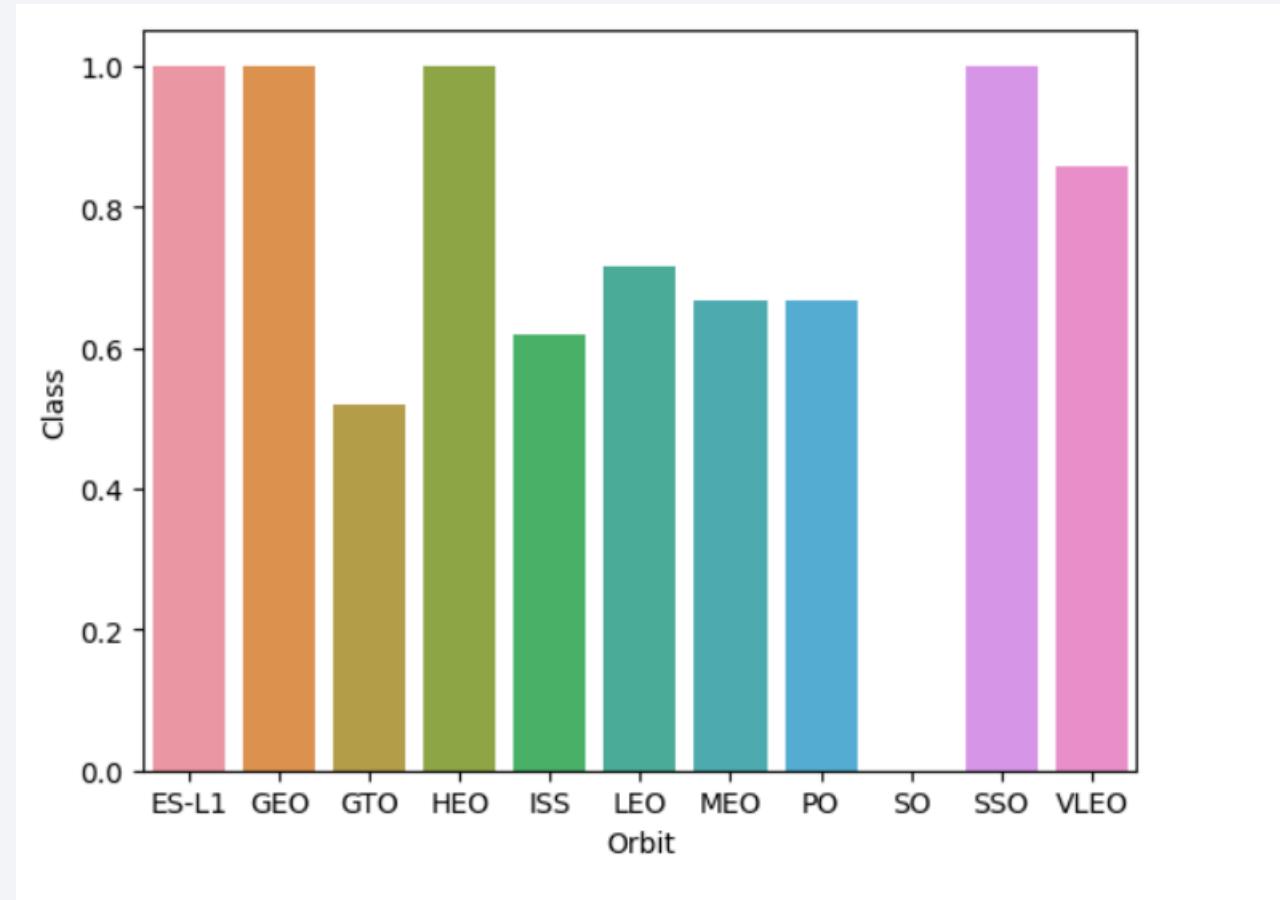


- The greater the payload mass for CCAFS SLC 40 the higher the success rate.

# Success Rate vs. Orbit Type

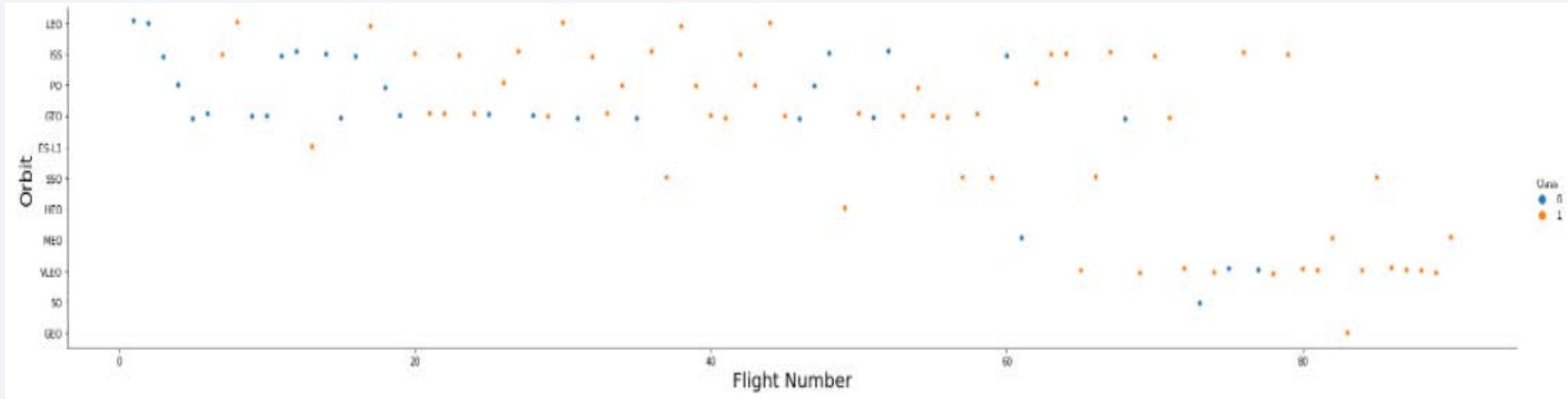
---

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



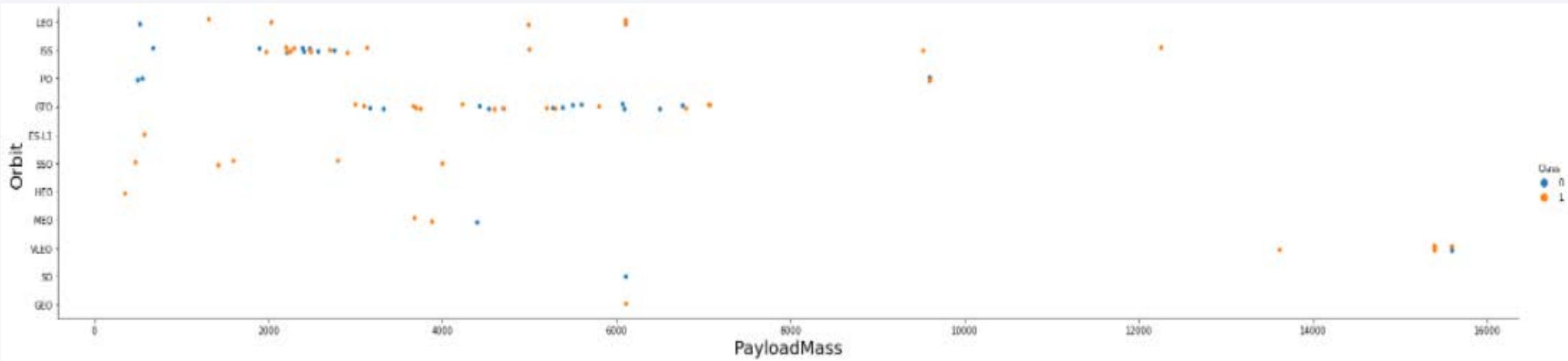
# Flight Number vs. Orbit Type

---



- 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.

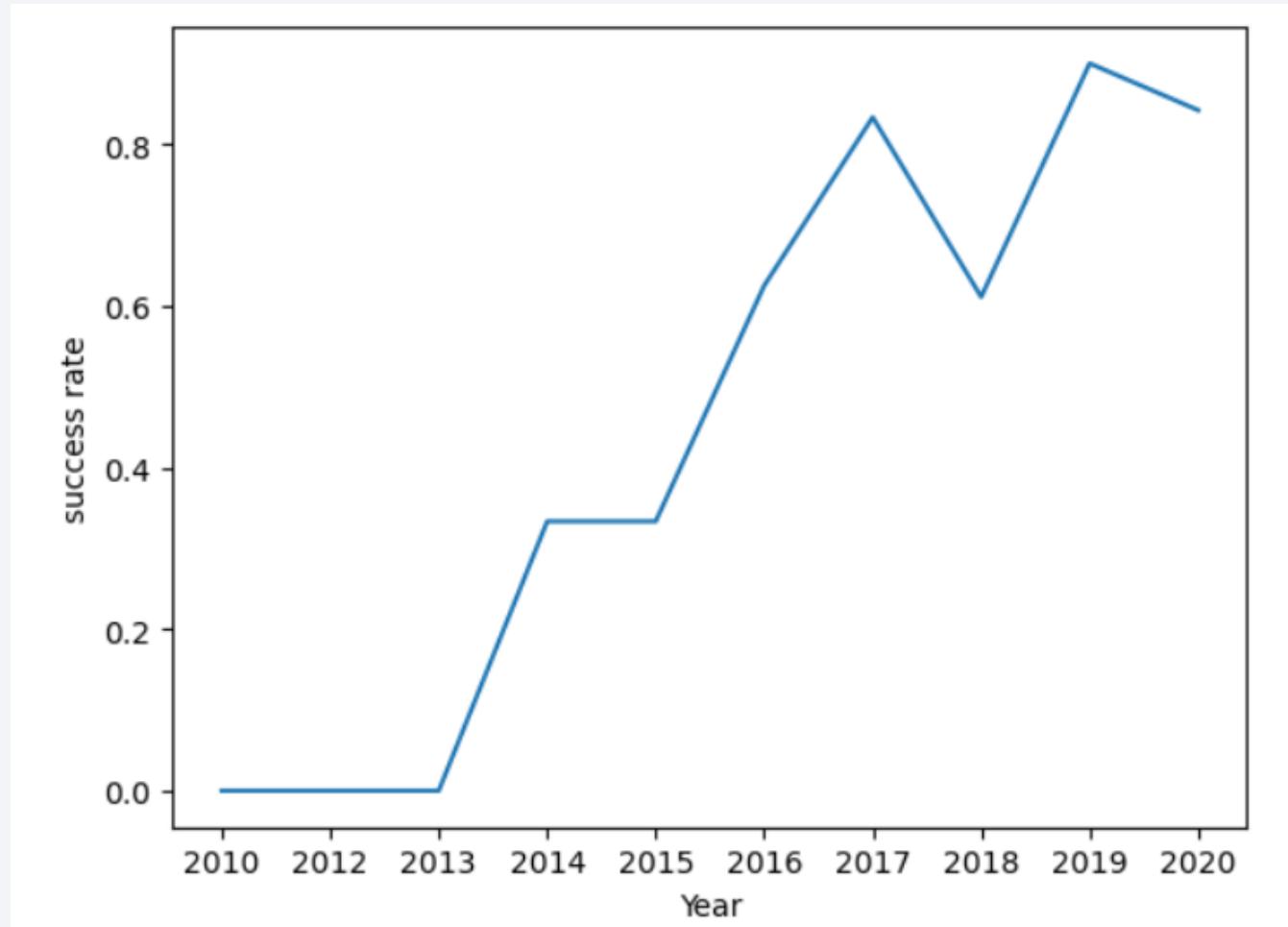
# Payload vs. Orbit Type



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

# Launch Success Yearly Trend

---



- 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.

```
%sql select distinct(LAUNCH_SITE) from SPACEXTBL;
```

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

```
Done.
```

## Launch\_Site

---

```
CCAFS LC-40
```

```
VAFB SLC-4E
```

```
KSC LC-39A
```

```
CCAFS SLC-40
```

# Launch Site Names Begin with 'CCA'

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

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

```
Done.
```

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

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

# Total Payload Mass

---

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

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

```
Done.
```

```
sum(PAYLOAD_MASS__KG_)
```

```
45596
```

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

# Average Payload Mass by F9 v1.1

---

```
%sq1 select avg(PAYLOAD_MASS__KG_) from SPACEEXTBL where Booster_Version = 'F9 v1.1';  
* sqlite:///my_data1.db  
Done.  
avg(PAYLOAD_MASS__KG_)  
2928.4
```

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

# First Successful Ground Landing Date

---

```
%sql select min(date) from SPACEXTBL where "LANDING _OUTCOME" = 'Success (ground pad)';
```

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

```
Done.
```

min(date)
2015-12-22

- We observed that the dates of the first successful landing outcome on ground pad was 22<sup>nd</sup> December 2015

## Successful Drone Ship Landing with Payload between 4000 and 6000

---

```
%sql select BOOSTER_VERSION from SPACEXTBL where "LANDING _OUTCOME" ='Success (drone ship)' and PAYLOAD_MASS__KG_ BETWEEN 4000 AND 6000;
```

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

### Booster\_Version

F9 FT B1022
F9 FT B1026
F9 FT B1021.2
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

```
%sql1 Select MISSION_OUTCOME, count(MISSION_OUTCOME) as count from SPACEXTBL GROUP BY MISSION_OUTCOME;
```

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

```
Done.
```

Mission_Outcome	count
Failure (in flight)	1
Success	98
Success	1
Success (payload status unclear)	1

- We used count to calculate the total Number of Successful and Failure

# Boosters Carried Maximum Payload

```
%sql select distinct BOOSTER_VERSION from SPACEXTBL where PAYLOAD_MASS_KG_ = (select MAX(PAYLOAD_MASS_KG_) FROM SPACEXTBL);
```

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

```
Done.
```

## 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
```

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

# 2015 Launch Records

---

```
%sql select "LANDING _OUTCOME", BOOSTER_VERSION, LAUNCH_SITE, substr(Date,7,4) as year from SPACEXTBL where year='2015' and "LANDING _OUTCOME" = 'Failure'
```

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

```
Done.
```

Landing_Outcome	Booster_Version	Launch_Site	year
Failure (drone ship)	F9 v1.1 B1012	CCAFS LC-40	2015
Failure (drone ship)	F9 v1.1 B1015	CCAFS LC-40	2015

- We used the **WHERE** clause to filter for failed landing outcomes in drone ship, their booster versions, and launch site names for year 2015

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

---

```
%sql select "LANDING _OUTCOME", count("LANDING _OUTCOME") as count from SPACEXTBL group by "LANDING _OUTCOME" having date between '04-06-2010' and '2017-03-20'
```

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

```
Done.
```

Landing _Outcome	count
No attempt	10
Success (drone ship)	6
Failure (drone ship)	5
Success (ground pad)	5
Controlled (ocean)	3
Uncontrolled (ocean)	2
Precalculated (drone ship)	1
Failure (parachute)	1

Landing _Outcome	count
No attempt	10
Success (drone ship)	6
Failure (drone ship)	5
Success (ground pad)	5
Controlled (ocean)	3
Uncontrolled (ocean)	2
Precalculated (drone ship)	1
Failure (parachute)	1

- 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 the dark void of space. City lights are visible as numerous small white and yellow dots, primarily concentrated in coastal and urban areas. The atmosphere appears as a thin blue layer above the clouds, which are depicted as dark, swirling patterns.

Section 3

# Launch Sites Proximities Analysis

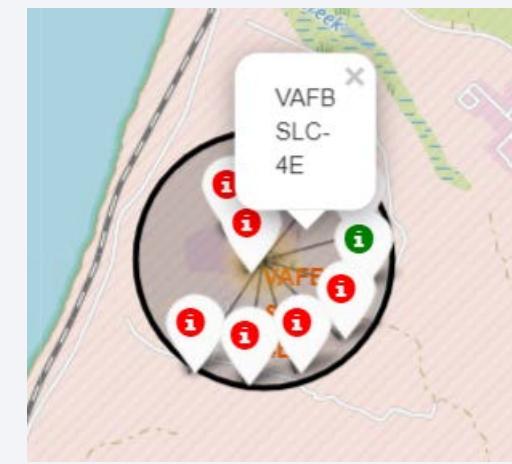
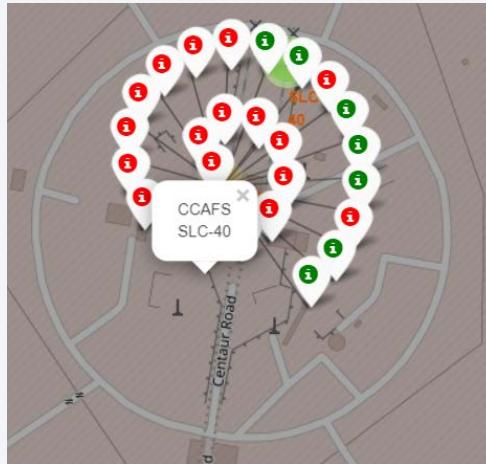
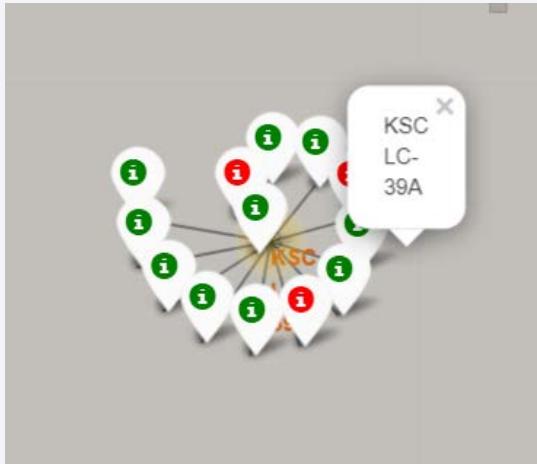
# All launch sites global map markers



- We can see that the SpaceX launch sites are in the USA coasts

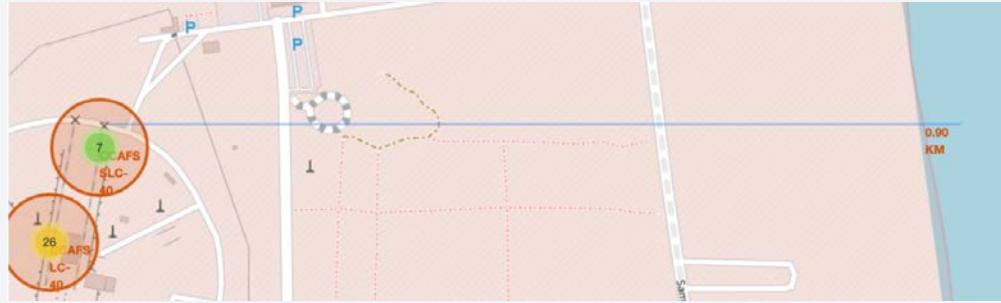
# Markers showing launch sites with color labels

---

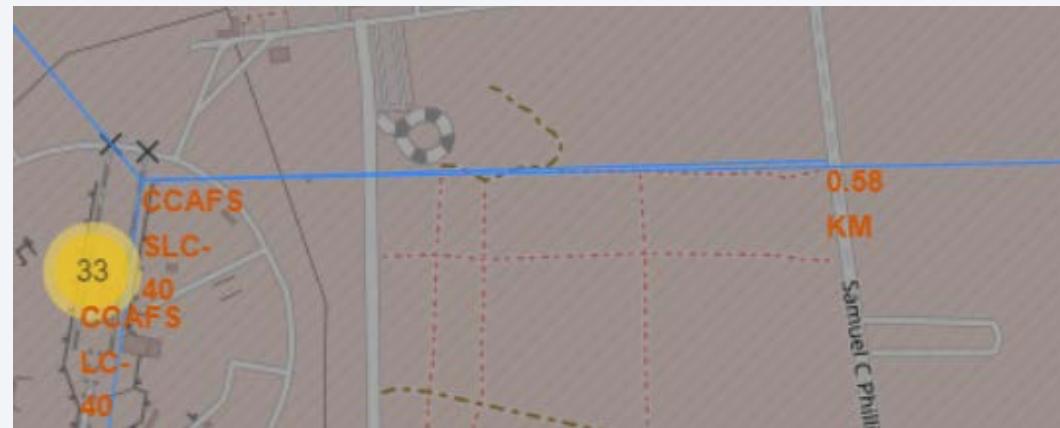
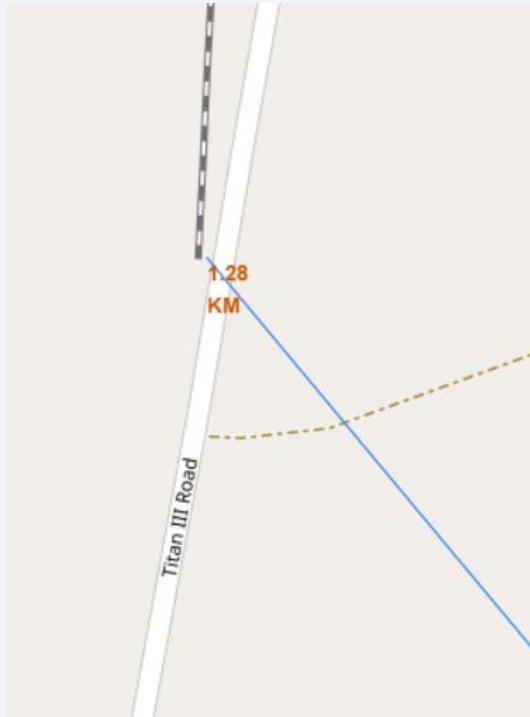
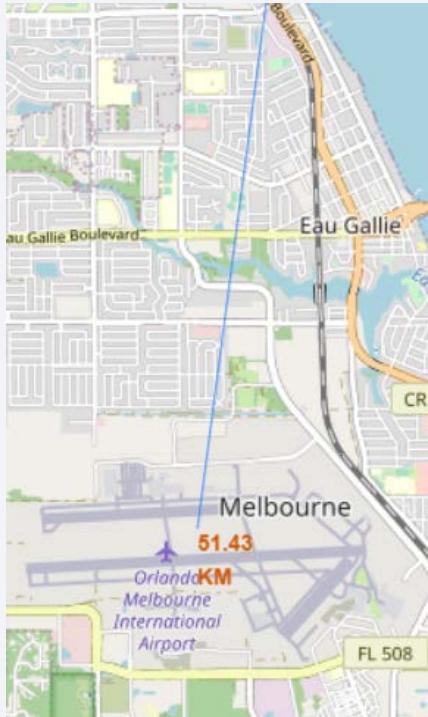


- Green Marker shows Success; Red Marker shows Failures

# Launch Site distance to landmarks

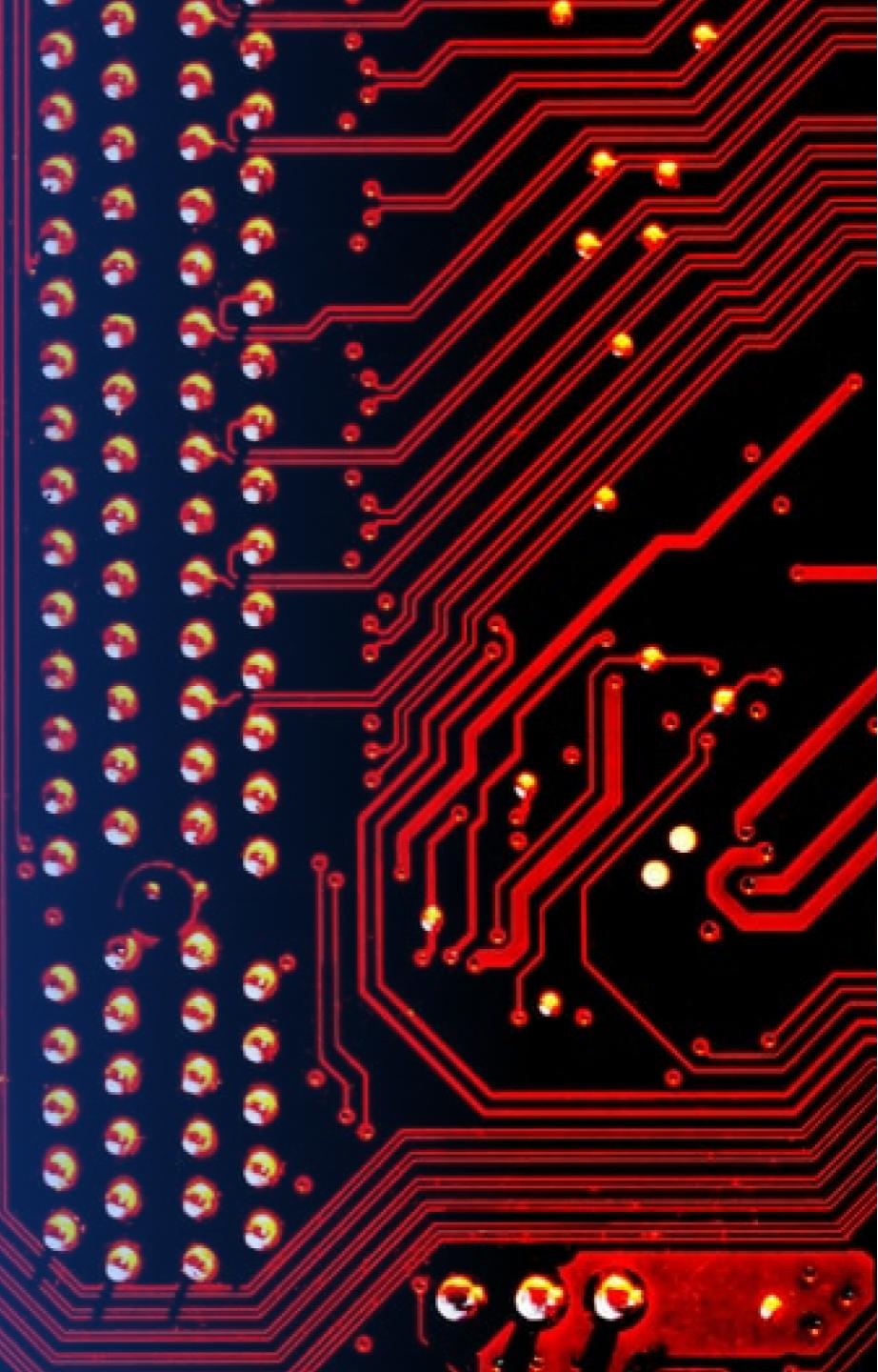


- Launch sites are in close proximity to coastline, and keep certain distance away from railways, highways and cities.



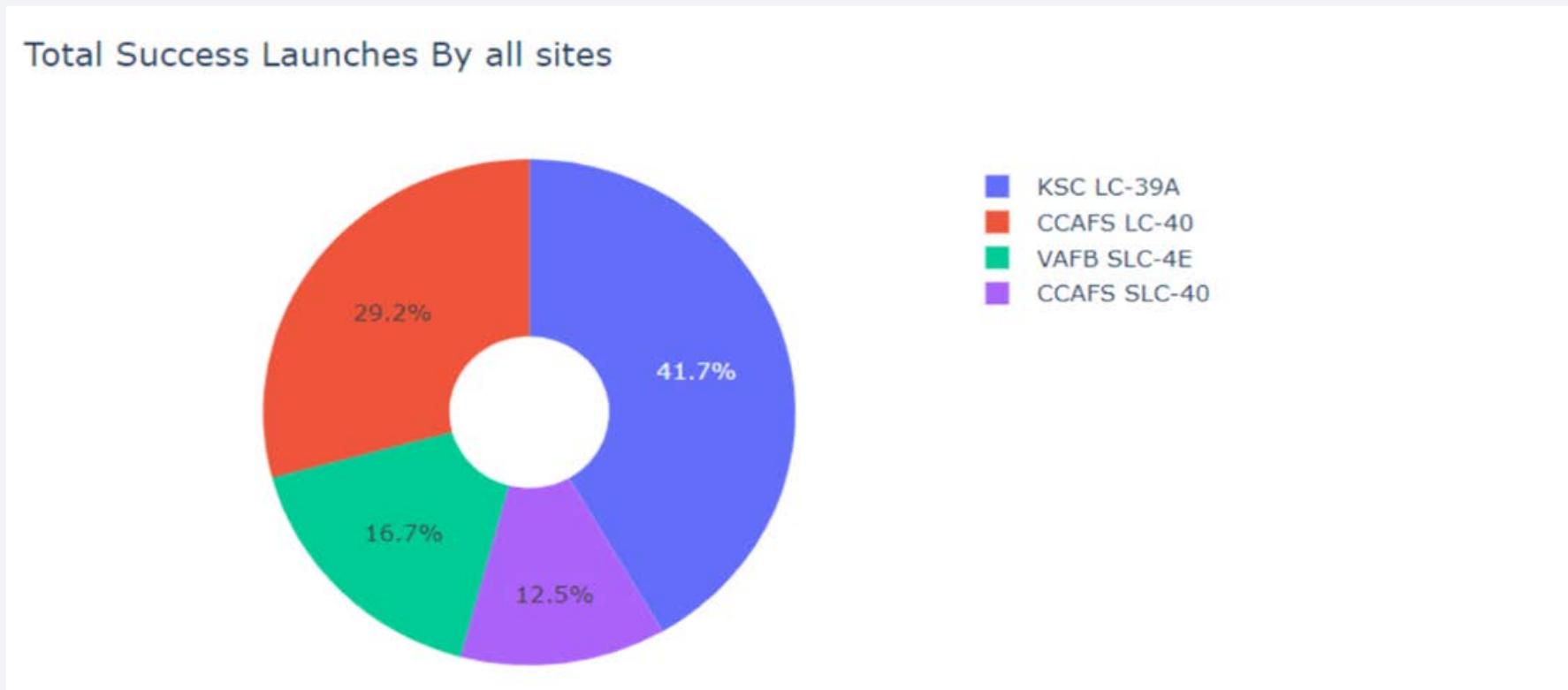
Section 4

# Build a Dashboard with Plotly Dash



## Pie chart showing the success percentage achieved by each launch site

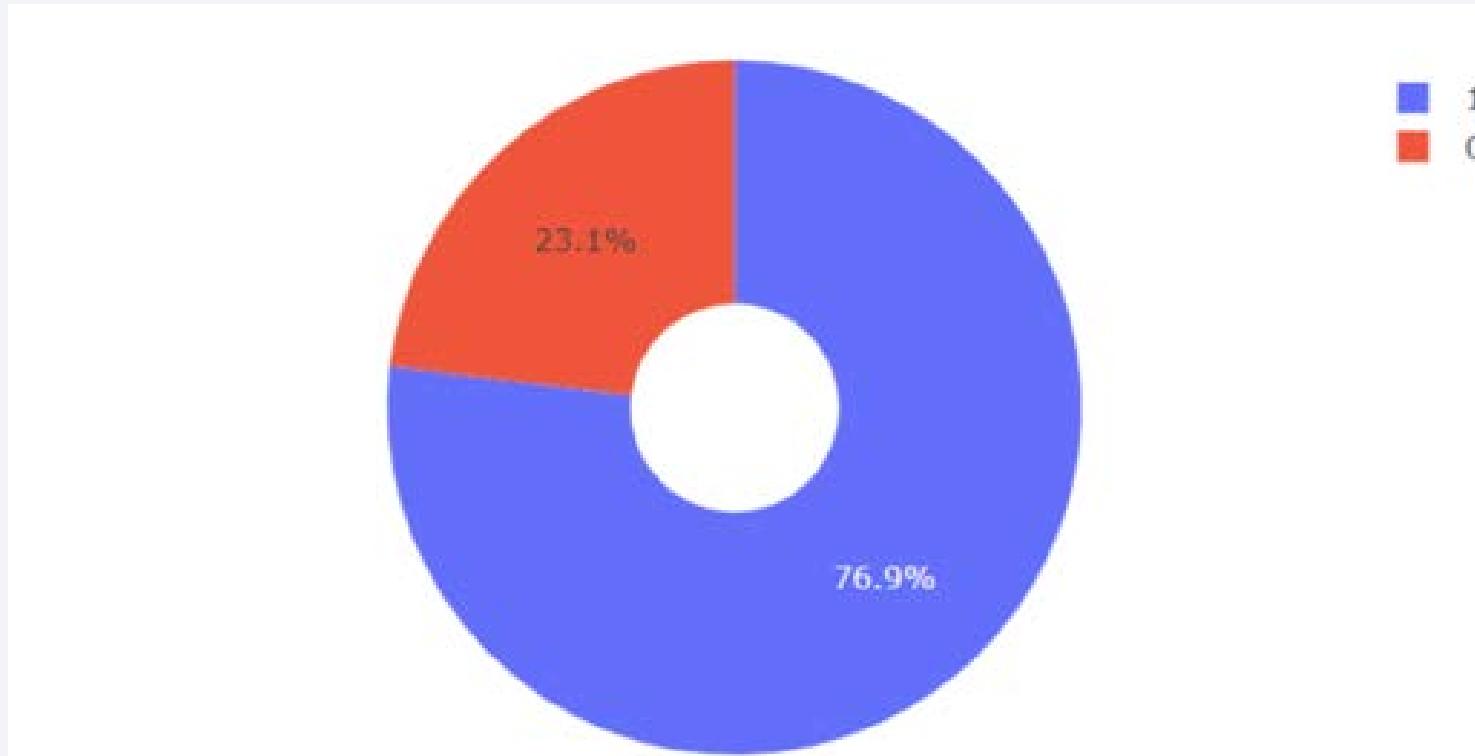
---



- KSC LC-39A had most successful launches from all the sites

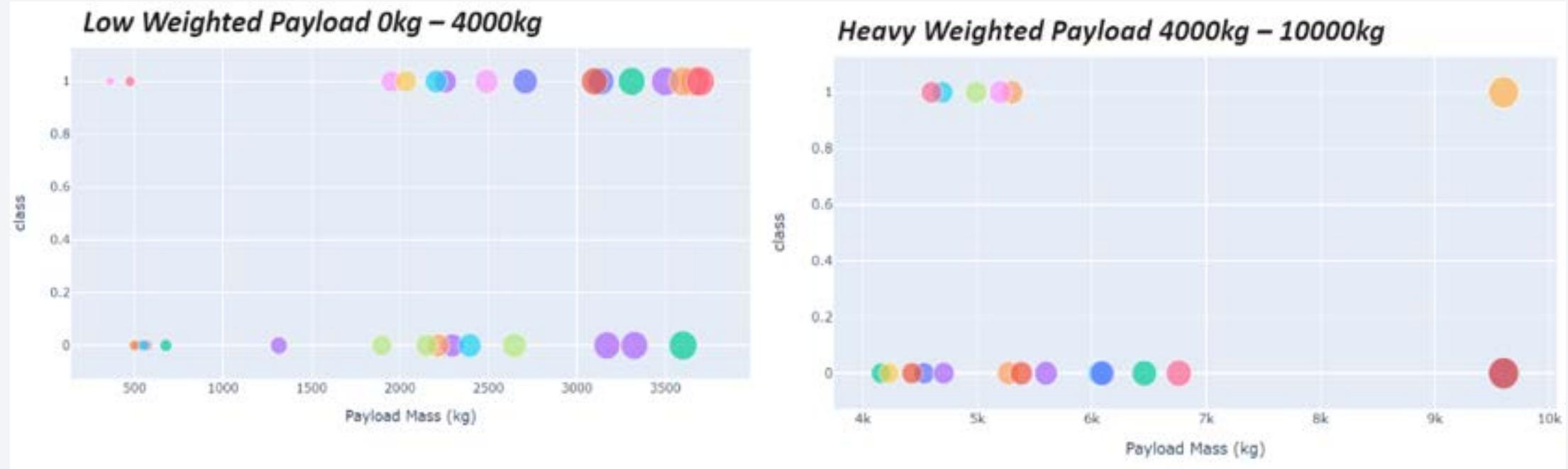
## Pie chart showing the Launch site with the highest launch success ratio

---



- KSC LC-39A achieved 76.9% success rate

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



- We can see the success rates for low weighted payloads is higher than heavy

Section 5

# Predictive Analysis (Classification)

# Classification Accuracy

---

```
print('Accuracy for LR method:', logreg_cv.score(X_test, Y_test))
print('Accuracy for SVM method:', svm_cv.score(X_test, Y_test))
print('Accuracy for Tree method:', tree_cv.score(X_test, Y_test))
print('Accuracy for KNN method:', knn_cv.score(X_test, Y_test))
```

Accuracy for LR method: 0.8333333333333334

Accuracy for SVM method: 0.8333333333333334

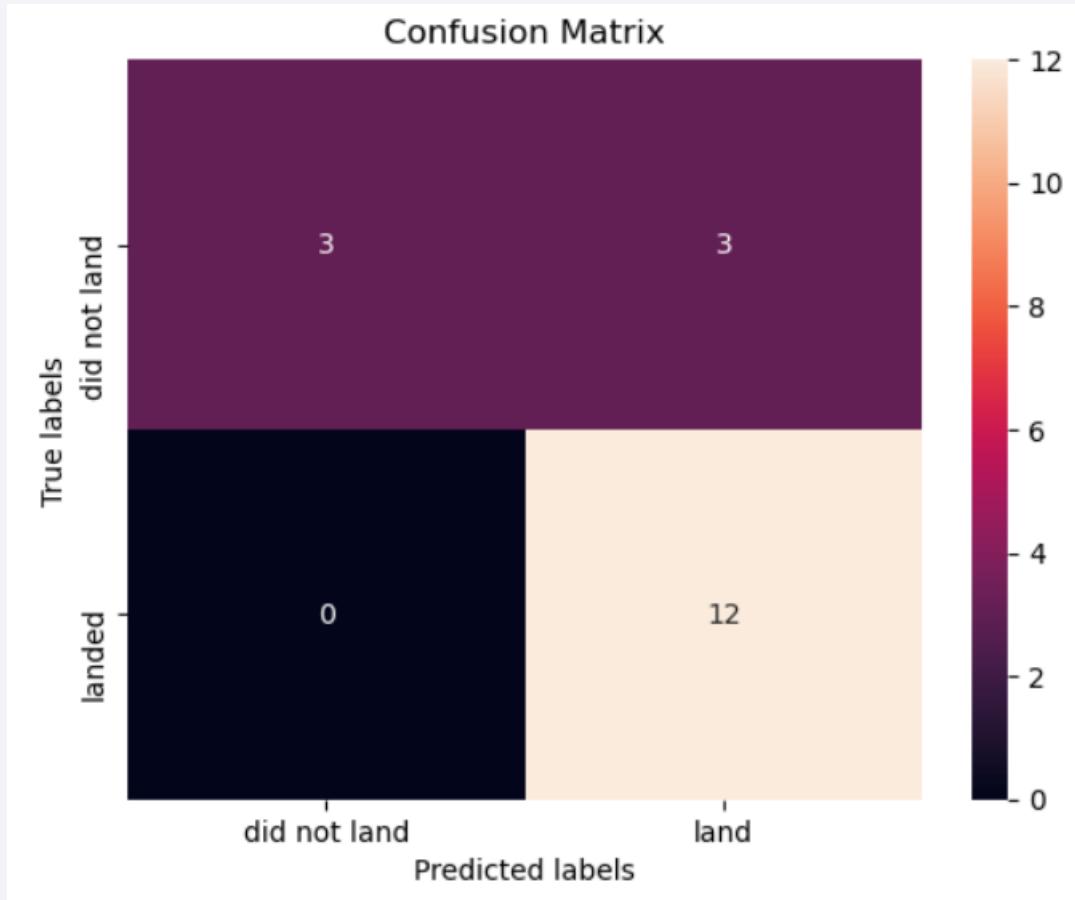
Accuracy for Tree method: 0.9444444444444444

Accuracy for KNN method: 0.8333333333333334

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

# Confusion Matrix

---



- The confusion matrix for the decision tree classifier shows that the classifier can distinguish between the different classes.
- The major problem is that unsuccessful landing marked as successful landing by the classifier (three events).

# Conclusions

---

We can conclude that:

- The larger the amount of flight, the greater the success rate at a launch site.
- Launch success rate had the rising trend in 2013 till 2020.
- Orbits ES-L1, GEO, HEO, SSO, VLEO had the most success rate.
- KSC LC-39A is the launch site which had the highest successful launch rate than others.
- The Decision tree classifier is the best machine learning algorithm for this task.(accuracy:94.4%)

Thank you!

