



IBM Developer  
SKILLS NETWORK

# Winning Space Race with Data Science

<Yuki Nakayama>  
<1st,MAR,2024>

in process



# Outline

---

**EXECUTIVE SUMMARY**

**INTRODUCTION**

**METHODOLOGY**

**RESULTS**

**CONCLUSION**

**APPENDIX**



# Executive Summary



## Summary of methodologies

1. Data Collection:
2. Data Preprocessing:
3. Exploratory Data Analysis (EDA):
4. Feature Engineering:
5. Model Selection:
6. Model Training:
7. Model Evaluation:

## Summary of all results

- Developed predictive models to forecast the success or failure of future SpaceX launches based on historical data, achieving a high level of accuracy and reliability. Utilized machine learning
- algorithms such as logistic regression and random forests to predict launch outcomes and assess the likelihood of mission success.

# 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.

## Problems you want to find answers

- What factors contribute to the success of landing?







Section 1

# Methodology

# Methodology

---

## Executive Summary

- Data collection methodology:
  - Accessing SpaceX data using the SpaceX API.
  - Retrieving information about launches, rockets, payloads, and other relevant data.
- Perform data wrangling
  - Cleaning and formatting the retrieved data. Handling missing values, duplicates, and outliers. Transforming data into a suitable format for analysis.
- Perform exploratory data analysis (EDA) using visualization and SQL
- Perform interactive visual analytics using Folium and Plotly Dash
- Perform predictive analysis using classification models
  - Create classification models with algorithms like logistic regression, decision trees, or random forests, fine-tune parameters, and evaluate accuracy using metrics such as accuracy and F1-score.

# Data Collection

---



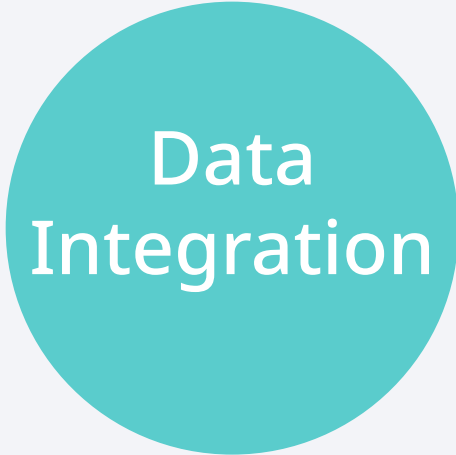
## Access Data Sources

Utilize APIs to access real-time data from SpaceX, ensuring up-to-date information on launches, rockets, and payloads.



## Data Cleaning and Preprocessing

Perform data cleaning to remove duplicates, handle missing values, and address inconsistencies in the dataset.



## Data Integration

Merge and integrate multiple datasets obtained from different sources to create a unified dataset for analysis.



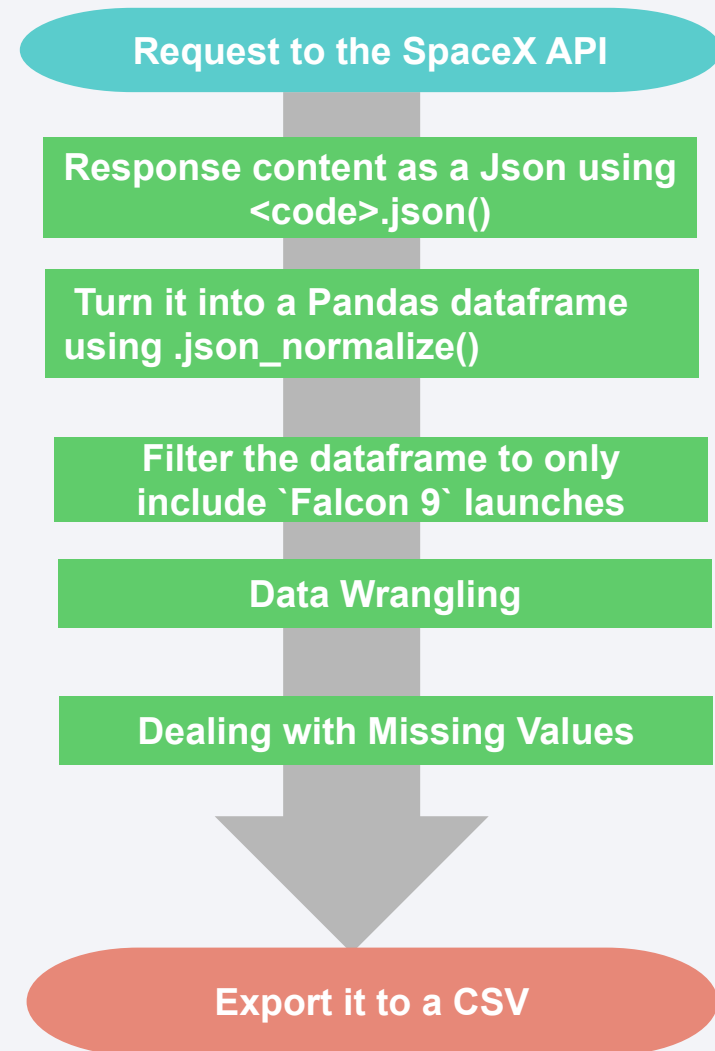
## Data Validation

Validate the integrity and accuracy of the collected data by cross-referencing with trusted sources and performing data quality checks.

# Data Collection – SpaceX API

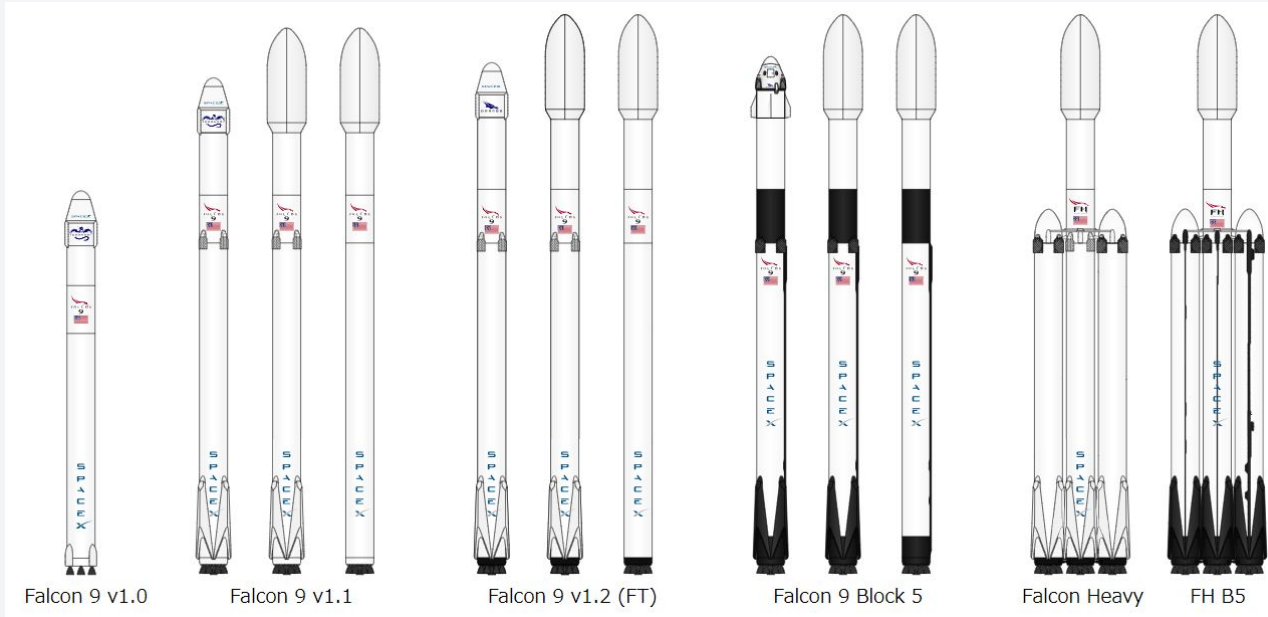
```
Now let's start requesting rocket launch data from SpaceX API with the following URL:  
[6]: spacex_url="https://api.spacexdata.com/v4/launches/past"  
[7]: response = requests.get(spacex_url)
```

[https://github.com/YukiG16/IBM-Data-Science-Certificate/blob/main/Course10\\_applied-data-science-capstone/Data-science-using-SpaceX-API/week1-1\\_jupyter-labs-spacex-data-collection-api.ipynb](https://github.com/YukiG16/IBM-Data-Science-Certificate/blob/main/Course10_applied-data-science-capstone/Data-science-using-SpaceX-API/week1-1_jupyter-labs-spacex-data-collection-api.ipynb)





# Data Collection - Scraping



[https://github.com/YukiG16/IBM-Data-Science-Certificate/blob/main/Course10\\_applied-data-science-capstone/Data-science-using-SpaceX-API/week1-2\\_jupyter-labs-webscraping.ipynb](https://github.com/YukiG16/IBM-Data-Science-Certificate/blob/main/Course10_applied-data-science-capstone/Data-science-using-SpaceX-API/week1-2_jupyter-labs-webscraping.ipynb)

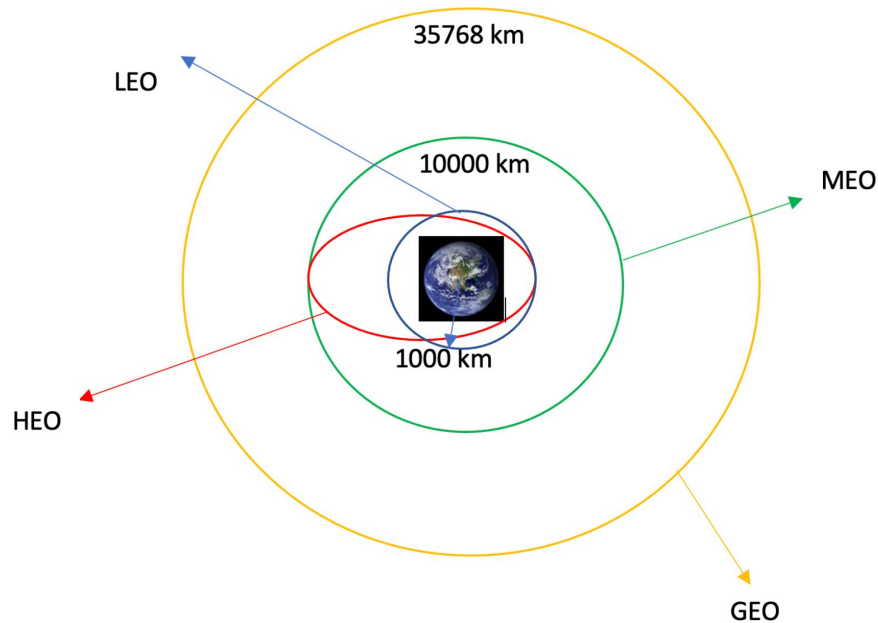
Request the Falcon9 Launch Wiki page from its URL

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

Extract all column/variable names from the HTML table head

Export it to a CSV

# Data Wrangling



Each launch aims to an dedicated orbit

[https://github.com/YukiG16/IBM-Data-Science-Certificate/blob/main/Course10\\_applied-data-science-capstone/Data-science-using-SpaceX-API/week1-3\\_labs-jupyter-spacex-Data%20wrangling.ipynb](https://github.com/YukiG16/IBM-Data-Science-Certificate/blob/main/Course10_applied-data-science-capstone/Data-science-using-SpaceX-API/week1-3_labs-jupyter-spacex-Data%20wrangling.ipynb)

Yuki Nakayama 1st,MAR,2024

Calculate the number of launches  
on each site

Calculate the number and  
occurrence of each orbit

Calculate the number and  
occurrence of mission outcome of  
the orbits

Create a landing outcome label  
from Outcome column

Export it to a CSV

# EDA with Data Visualization

---

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 SpaceX for a rocket launch.

This dataset includes a record for each payload carried during a SpaceX mission into outer space.

[https://github.com/YukiG16/IBM-Data-Science-Certificate/blob/main/Course10\\_applied-data-science-capstone/Data-science-using-SpaceX-API/week2-1\\_jupyter-labs-eda-sql-coursera\\_sqlite.ipynb](https://github.com/YukiG16/IBM-Data-Science-Certificate/blob/main/Course10_applied-data-science-capstone/Data-science-using-SpaceX-API/week2-1_jupyter-labs-eda-sql-coursera_sqlite.ipynb)

# EDA with SQL

---

- Using bullet point format, summarize the SQL queries you performed

[https://github.com/YukiG16/IBM-Data-Science-Certificate/blob/main/Course10\\_applied-data-science-capstone/Data-science-using-SpaceX-API/week2-1\\_jupyter-labs-eda-sql-coursera\\_sqlite.ipynb](https://github.com/YukiG16/IBM-Data-Science-Certificate/blob/main/Course10_applied-data-science-capstone/Data-science-using-SpaceX-API/week2-1_jupyter-labs-eda-sql-coursera_sqlite.ipynb)



# Build an Interactive Map with Folium

---

- Summarize what map objects such as markers, circles, lines, etc. you created and added to a folium map
- Explain why you added those objects

in process

[https://github.com/YukiG16/IBM-Data-Science-Certificate/blob/main/Course10\\_applied-data-science-capstone/Data-science-using-SpaceX-API/week3-1\\_lab\\_jupyter\\_launch\\_site\\_location.jupyterlite.ipynb](https://github.com/YukiG16/IBM-Data-Science-Certificate/blob/main/Course10_applied-data-science-capstone/Data-science-using-SpaceX-API/week3-1_lab_jupyter_launch_site_location.jupyterlite.ipynb)

# Build a Dashboard with Plotly Dash

---

- Summarize what plots/graphs and interactions you have added to a dashboard
- Explain why you added those plots and interactions
- Add the GitHub URL of your completed Plotly Dash lab, as an external reference and peer-review purpose

**in process**

# Predictive Analysis (Classification)

---

- Summarize how you built, evaluated, improved, and found the best performing classification model
- You need present your model development process using key phrases and flowchart
- Add the GitHub URL of your completed predictive analysis lab, as an external reference and peer-review purpose

in process

# Results

---

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

**in process**



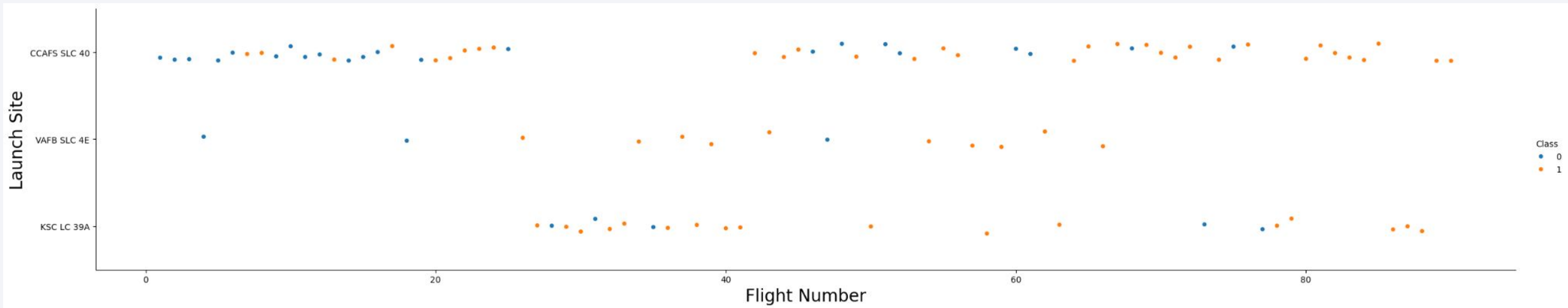


Section 2

# Insights drawn from EDA



# Flight Number vs. Launch Site



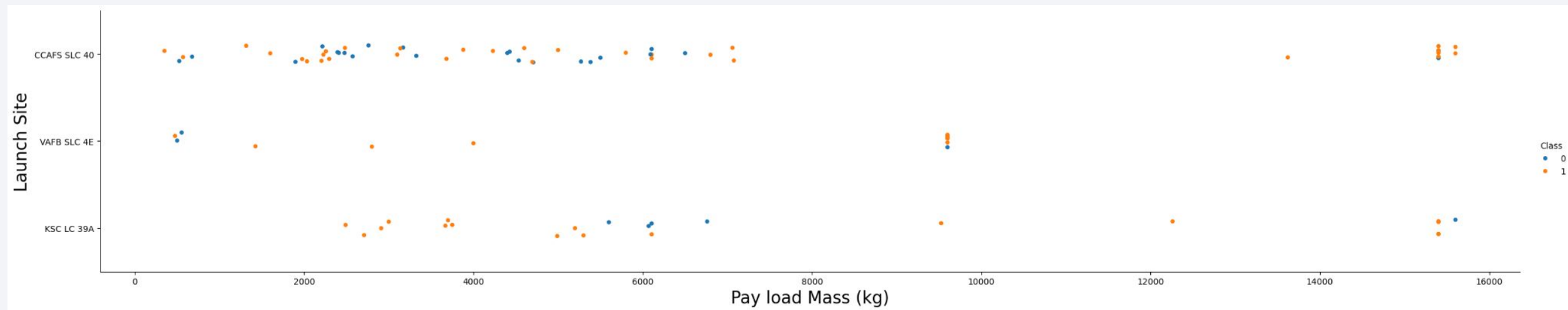
●:Success

●:failed

As the flight number increases, the first stage is more likely to land successfully.

This trend seems to be present across all launch sites.

# Payload vs. Launch Site

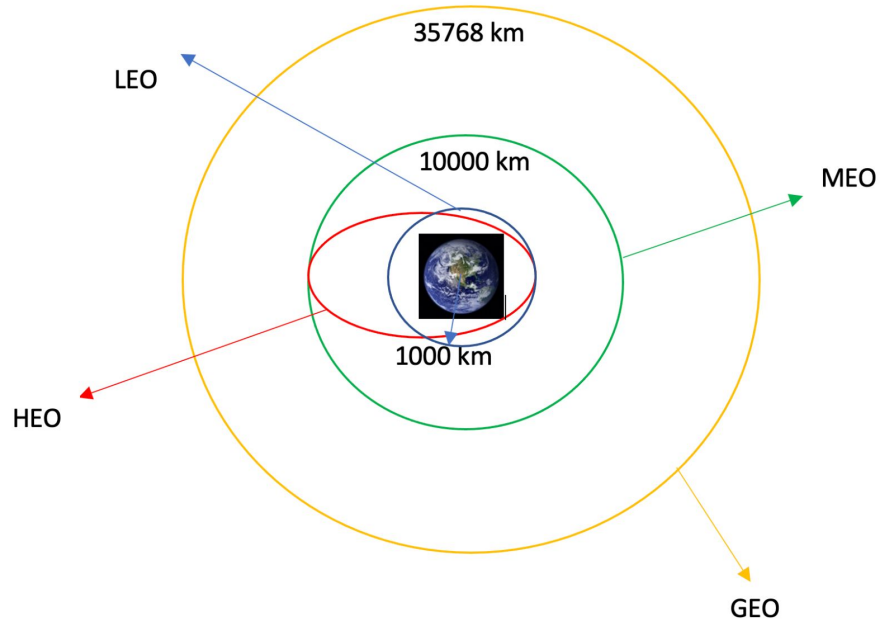


●:Success

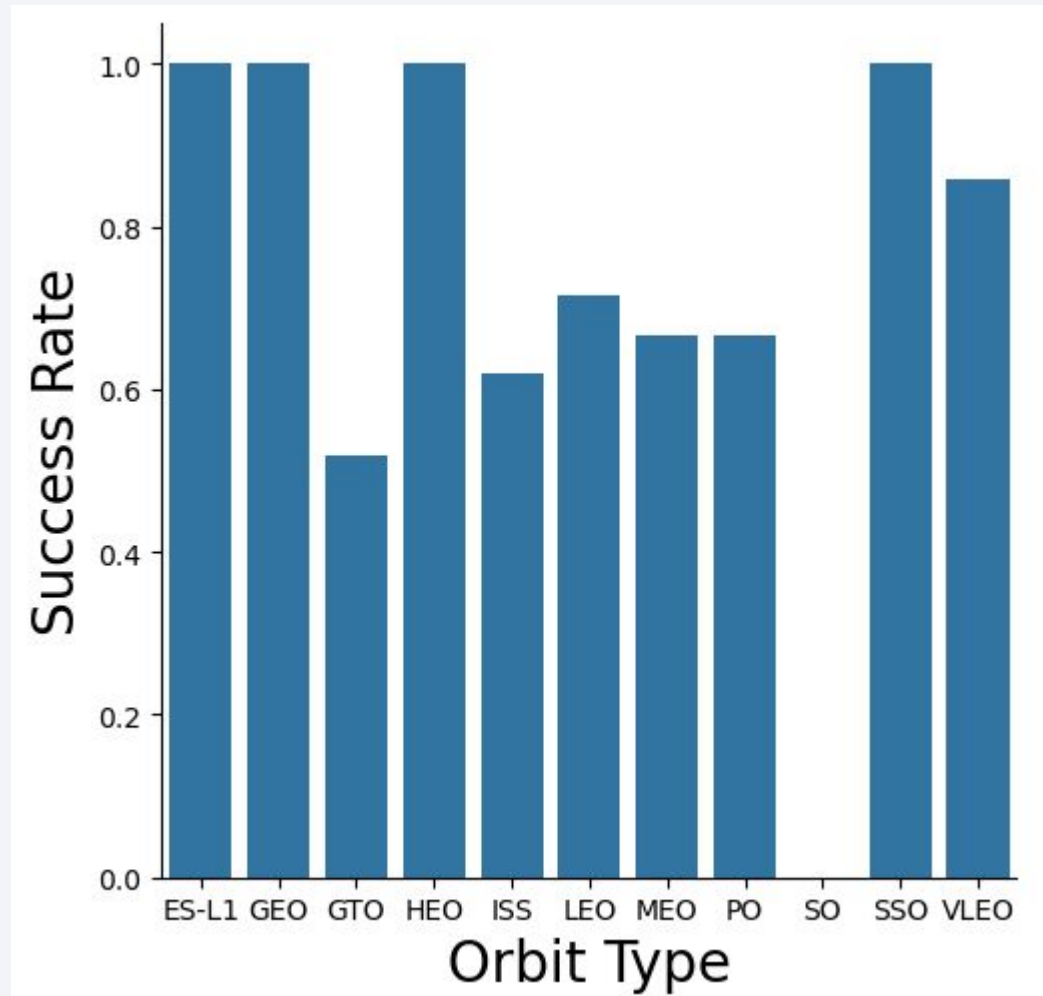
●:failed

It appears that as the payload increases, so does the flight number.

# Success Rate vs. Orbit Type

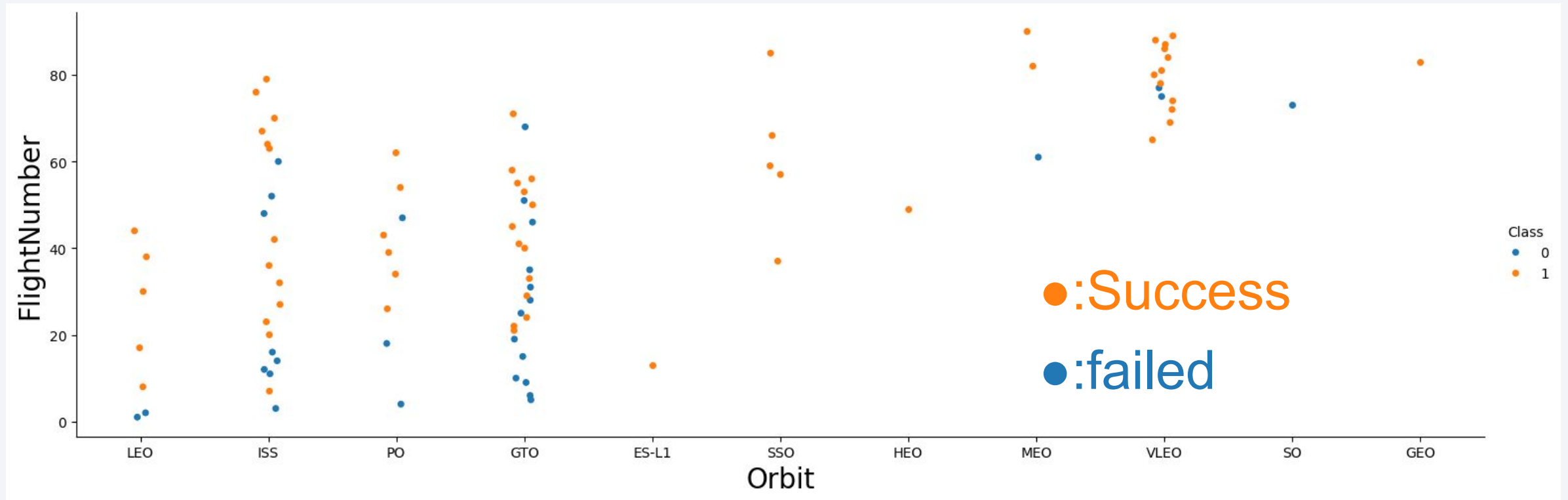


ES-L1, GEO, HEO, and SSO have a high success rate.



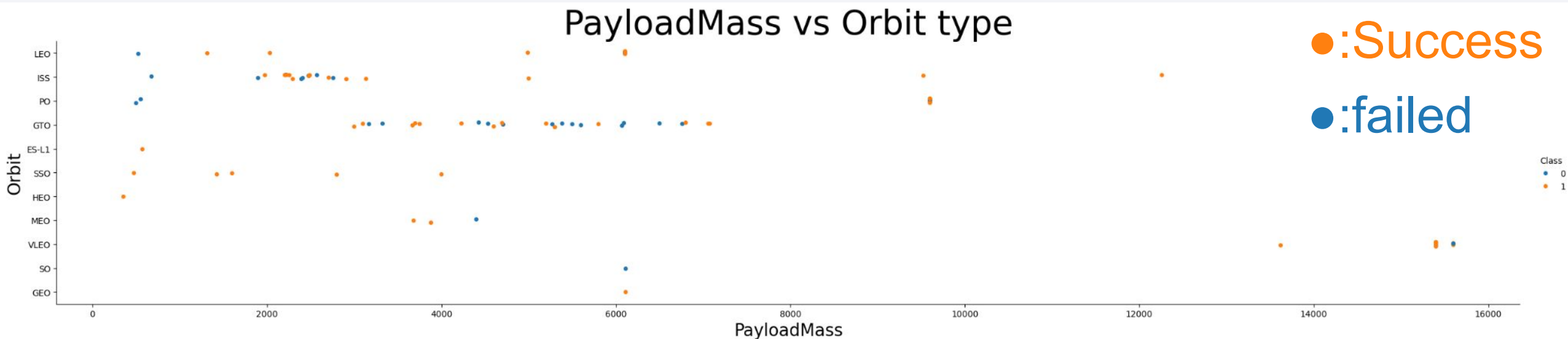


# Flight Number vs. Orbit Type



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



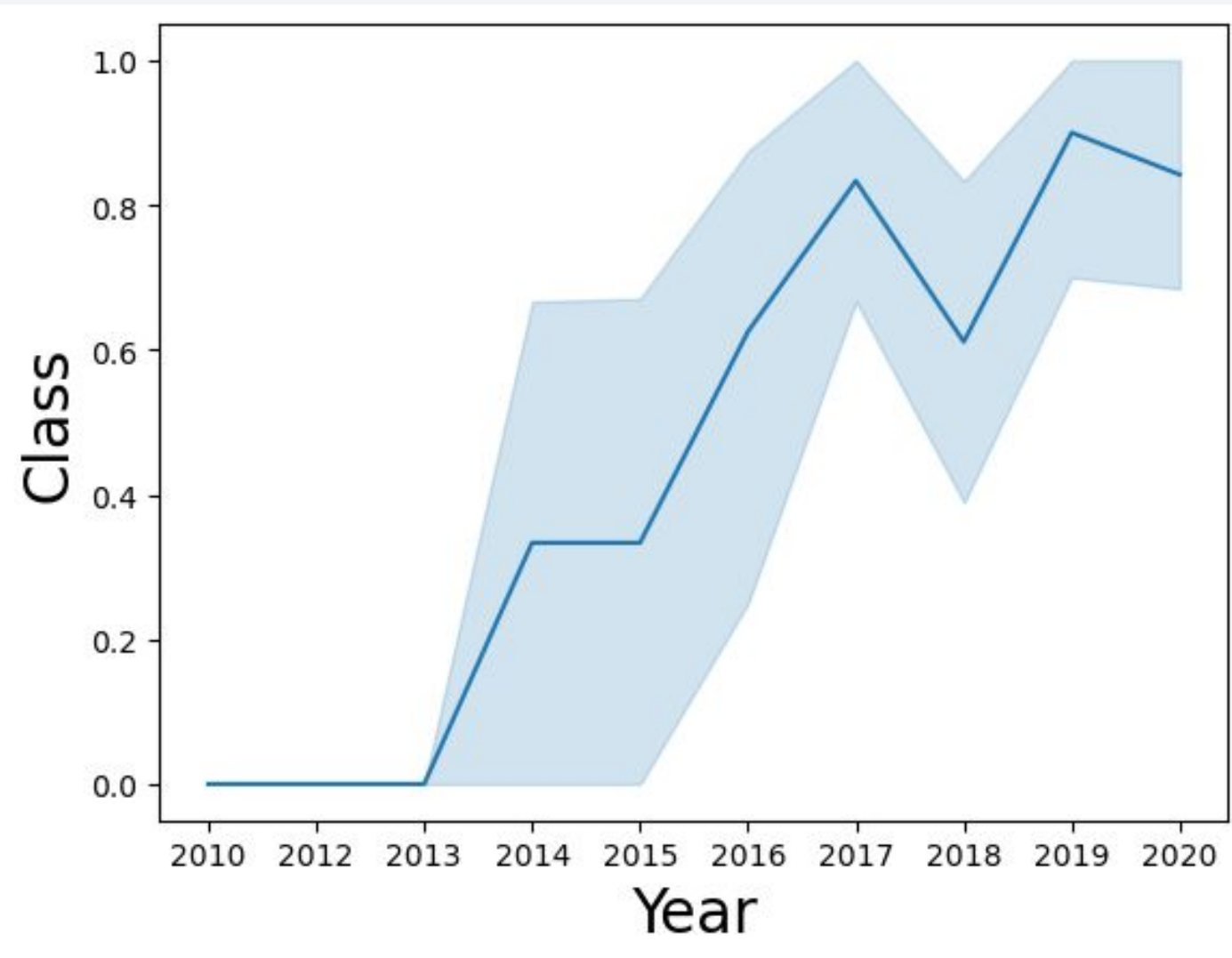
For heavy payloads, it appears that the successful landing rate or positive landing rate is higher for Polar, LEO, and ISS.

However for GTO we cannot distinguish this well as both positive landing rate and negative landing(unsuccesful mission) are both there here.

# Launch Success Yearly Trend

---

It can be noted that the success rate has shown a consistent increase from 2013 to 2020.



# All Launch Site Names

---

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

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

**SQL QUERY WITH A SHORT EXPLANATION**

Launch\_Site

CCAFS LC-40

VAFB SLC-4E

KSC LC-39A

CCAFS SLC-40



# Launch Site Names Begin with 'CCA'

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

```
%sql select * from SPACEXTBL 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	7:44:00	F9 v1.0 B0005	CCAFS LC-40	Dragon demo flight C2	525	LEO (ISS)	NASA (COTS)	Success	No attempt
2012-10-08	0:35:00	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

# Total Payload Mass

---

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

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

## SQL QUERY WITH A SHORT EXPLANATION

sum(PAYLOAD_MASS_KG_)
-----------------------

45596
-------

# Average Payload Mass by F9 v1.1

---

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

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

**SQL QUERY WITH A SHORT EXPLANATION**

avg(PAYLOAD\_MASS\_KG\_)

---

2928.4

# First Successful Ground Landing Date

---

List the date when the first succesful landing outcome in ground pad was acheived.

```
%sql select min(DATE) from SPACEXTBL where LANDING_OUTCOME = 'Success (ground pad)'
```

**SQL QUERY WITH A SHORT EXPLANATION**

**min(DATE)**

**2015-12-22**

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

---

List the names of the boosters which have success in drone ship and have payload mass greater than 4000 but less than 6000

```
%sql select BOOSTER_VERSION from SPACEXTBL where LANDING_OUTCOME = 'Success (drone ship)' and PAYLOAD_MASS_KG_ between 4000 and 6000
```

### SQL QUERY WITH A SHORT EXPLANATION

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

# Total Number of Successful and Failure Mission Outcomes

---

List the total number of successful and failure mission outcomes

```
%sql \  
SELECT \  
    Landing_Outcome, \  
    COUNT(*) AS Outcome_Count \  
FROM \  
    SPACEXTBL \  
WHERE \  
    Landing_Outcome IN ('Success', 'Failure') \  
GROUP BY \  
    Landing_Outcome;
```

**SQL QUERY WITH A SHORT EXPLANATION**

Landing_Outcome	Outcome_Count
Failure	3
Success	38



# Boosters Carried Maximum Payload

List the names of the booster\_versions which have carried the maximum payload mass. Use a subquery

```
%sql SELECT DISTINCT Booster_Version \  
FROM SPACEXTBL \  
WHERE PAYLOAD_MASS_KG_ = (SELECT MAX(PAYLOAD_MASS_KG_) FROM SPACEXTBL LIMIT 1);
```

## SQL QUERY WITH A SHORT EXPLANATION

### 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

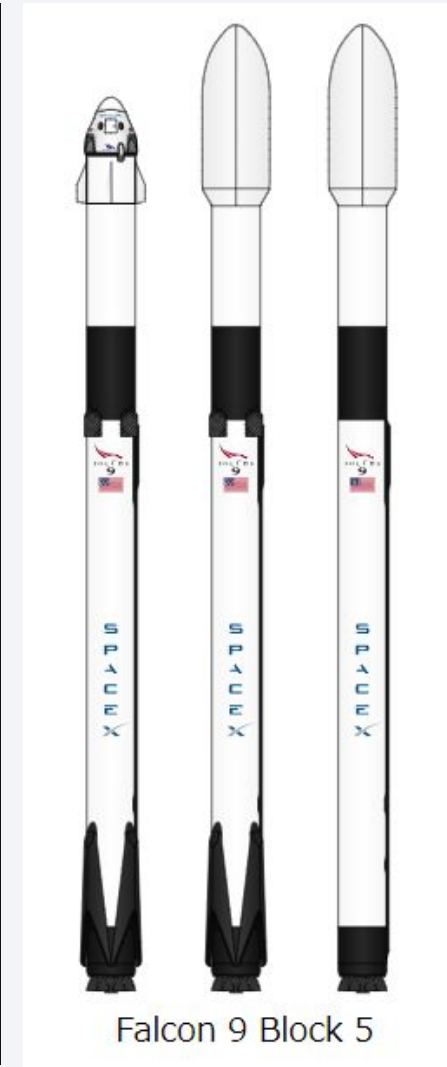
F9 B5 B1060.2

F9 B5 B1058.3

F9 B5 B1051.6

F9 B5 B1060.3

F9 B5 B1049.7



Falcon 9 Block 5

# 2015 Launch Records

Month	Booster_Version	Launch_Site	Landing_Outcome
01	F9 v1.1 B1012	CCAFS LC-40	Failure (drone ship)
04	F9 v1.1 B1015	CCAFS LC-40	Failure (drone ship)

List the records which will display the month names, failure landing\_outcomes in drone ship ,booster versions, launch\_site for the months in year 2015.

Note: SQLite does not support monthnames. So you need to use substr(Date, 6,2) as month to get the months and substr(Date,0,5)='2015' for year.

```
%sql \  
SELECT \  
    substr(Date, 6, 2) AS Month, \  
    Booster_Version, \  
    Launch_Site, \  
    Landing_Outcome \  
FROM \  
    SPACEXTBL \  
WHERE \  
    substr(Date, 0, 5) = '2015' \  
    AND Landing_Outcome LIKE 'Failure (drone ship)%';
```

## SQL QUERY WITH A SHORT EXPLANATION

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

Rank the count of landing outcomes (such as Failure (drone ship) or Success (ground pad)) between the date 2010-06-04 and 2017-03-20, in descending order.

```
%sql \  
SELECT \  
    Landing_Outcome, \  
    COUNT(*) AS Outcome_Count \  
FROM \  
    SPACEXTBL \  
WHERE \  
    Date BETWEEN '2010-06-04' AND '2017-03-20' \  
GROUP BY \  
    Landing_Outcome \  
ORDER BY \  
    Outcome_Count DESC;
```

## SQL QUERY WITH A SHORT EXPLANATION

Landing_Outcome	Outcome_Count
No attempt	10
Success (drone ship)	5
Failure (drone ship)	5
Success (ground pad)	3
Controlled (ocean)	3
Uncontrolled (ocean)	2
Failure (parachute)	2
Precluded (drone ship)	1



Section 3

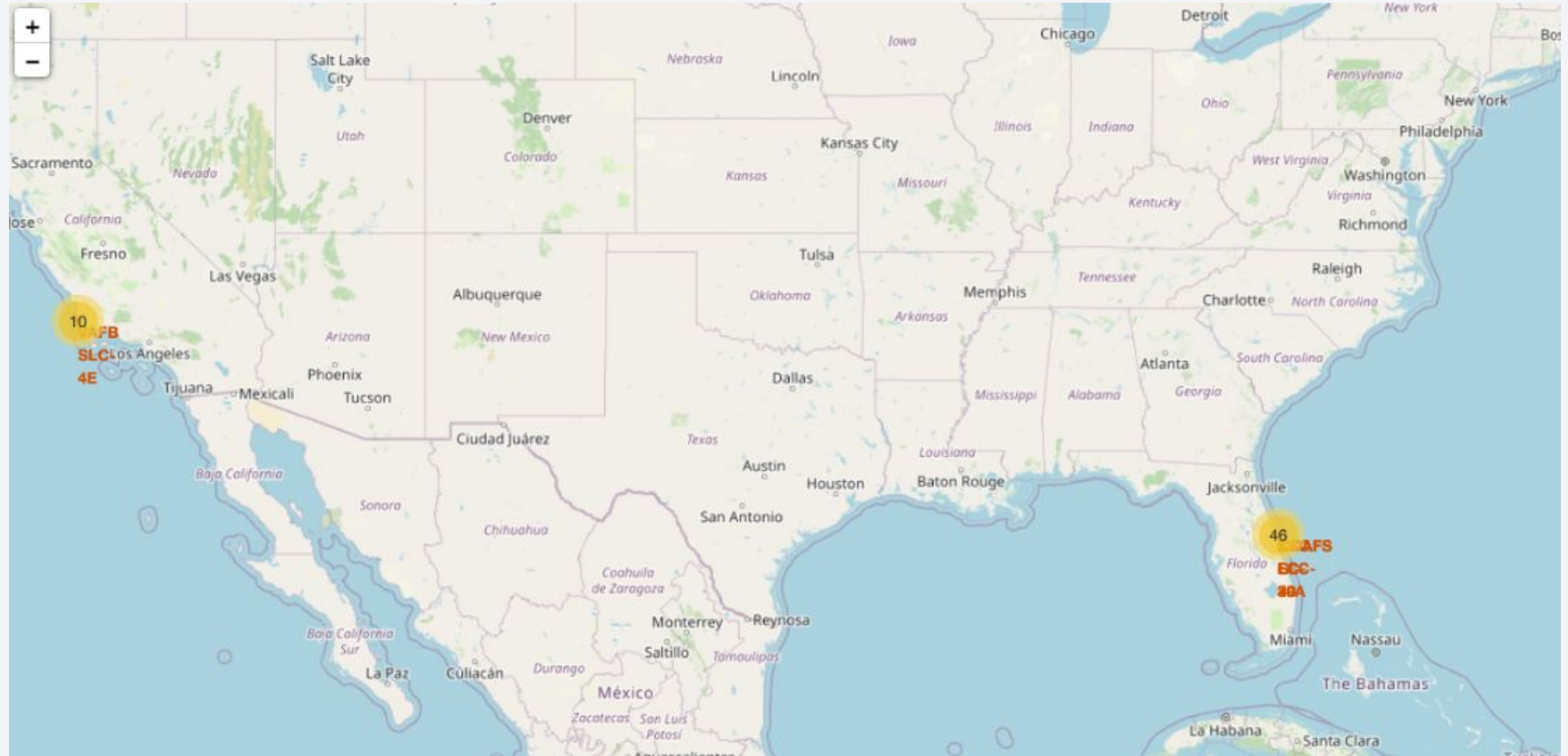
# Launch Sites Proximities Analysis



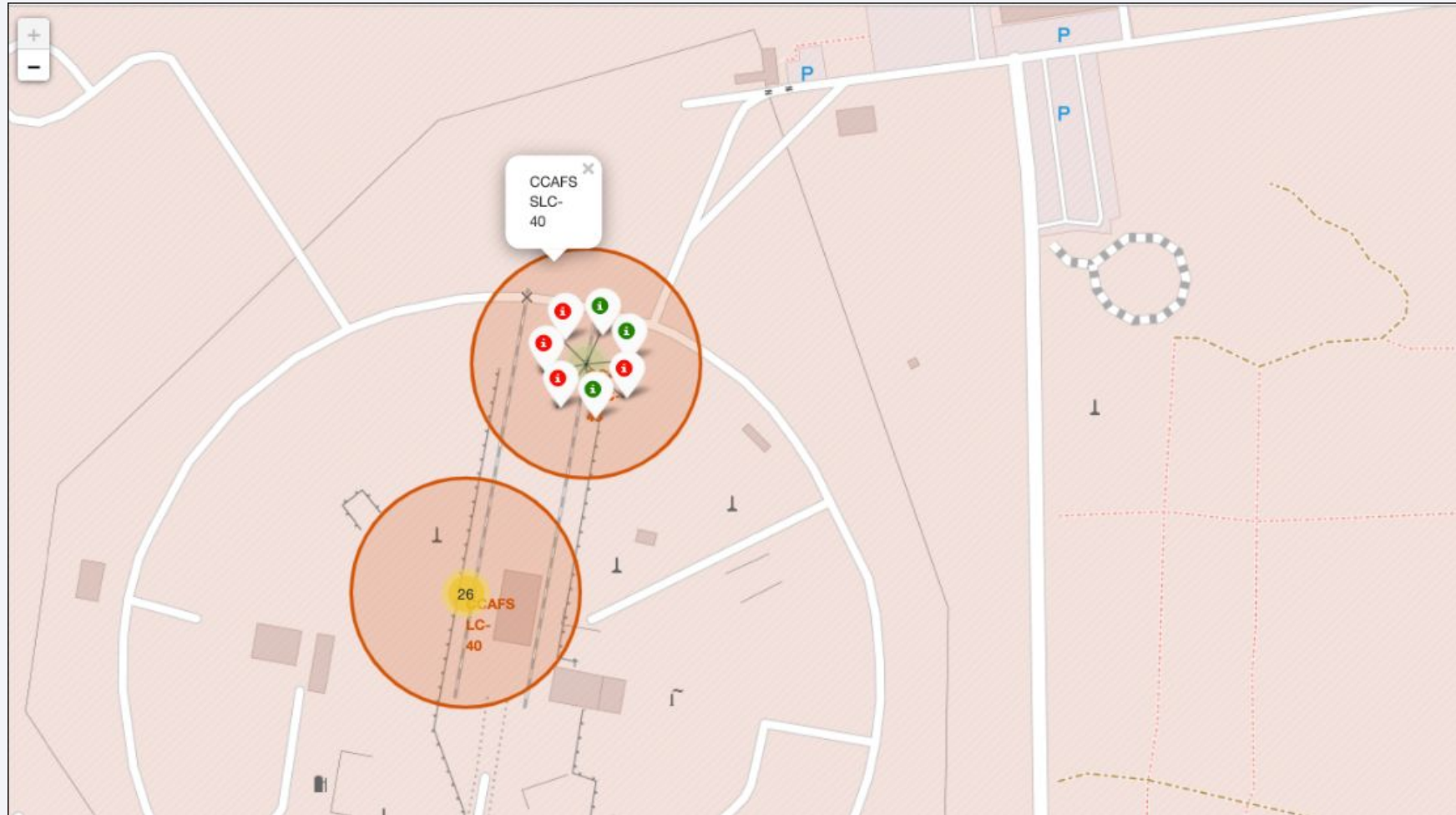
# Markers all launch sites on a map

**All launch sites are typically located near the Equator** due to the benefits of Earth's maximum rotational speed, which results in reduced fuel needs and lower costs for reaching orbit.

**All launch sites are strategically situated in close proximity to the coast** to ensure safe paths over water, thereby minimizing risks to populated areas during launch failures. In addition, it is worth noting that these systems can aid in the effective removal of rocket debris, while also providing logistical benefits for both operations and support.



## an overview of the outcomes of launches at each site on the map





# Distances between a launch site to its proximities

**Launch sites are in close proximity to railways :**  
Efficient transportation of heavy equipment.

**Launch sites are in close proximity to highways:**  
Easy access for personnel and support vehicles.

**Launch sites are in close proximity to coastline:**  
Clear trajectory over water, minimal risks to populated areas.

**launch sites keep certain distance away from cities:**  
Minimize risks to human life and property.

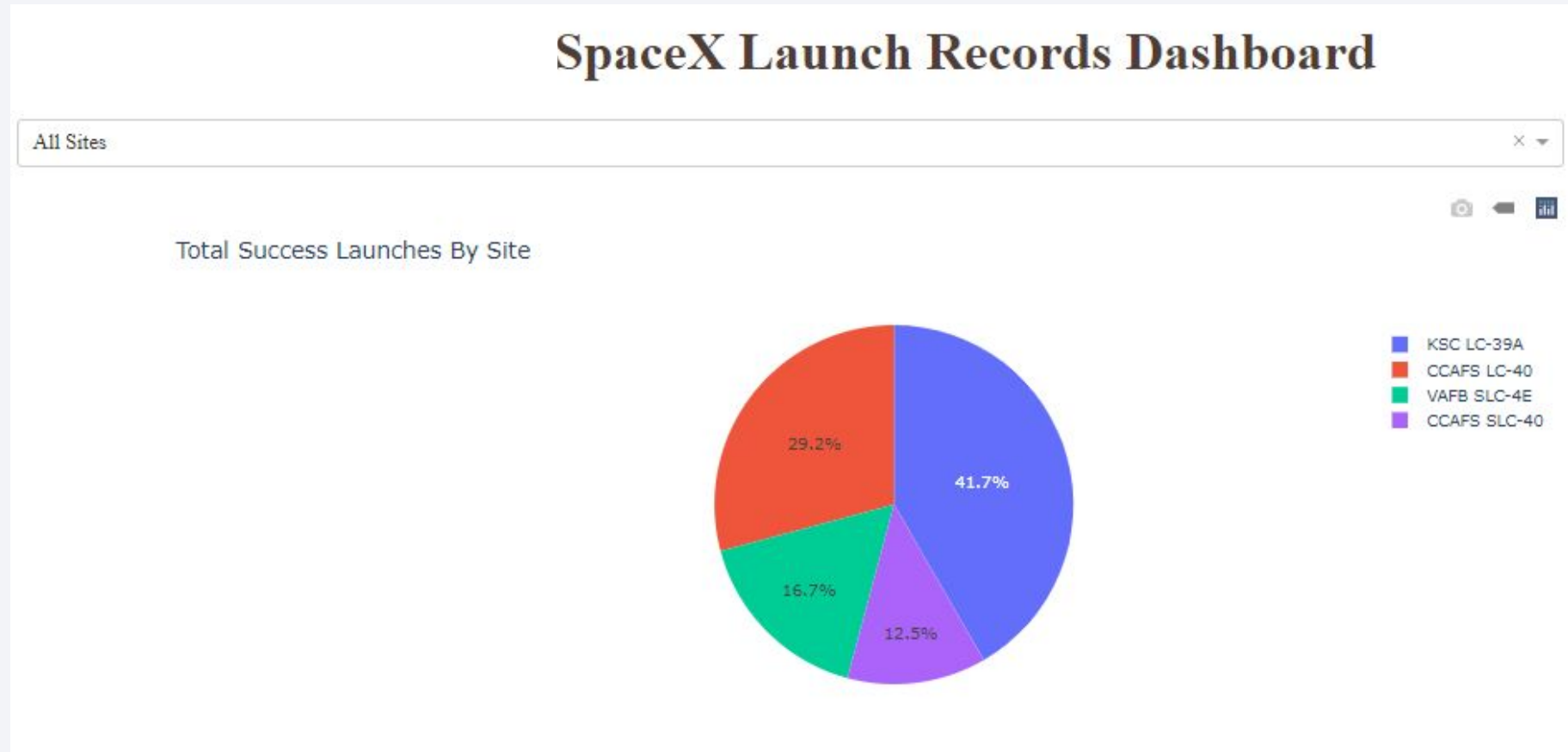




Section 4

# Build a Dashboard with Plotly Dash

# Total Success Launches By Site





# <Dashboard Screenshot 2>

---

- Replace <Dashboard screenshot 2> title with an appropriate title
- Show the screenshot of the piechart for the launch site with highest launch success ratio
- Explain the important elements and findings on the screenshot

## <Dashboard Screenshot 3>

---

- Replace <Dashboard screenshot 3> title with an appropriate title
- Show screenshots of Payload vs. Launch Outcome scatter plot for all sites, with different payload selected in the range slider
- Explain the important elements and findings on the screenshot, such as which payload range or booster version have the largest success rate, etc.



Section 5

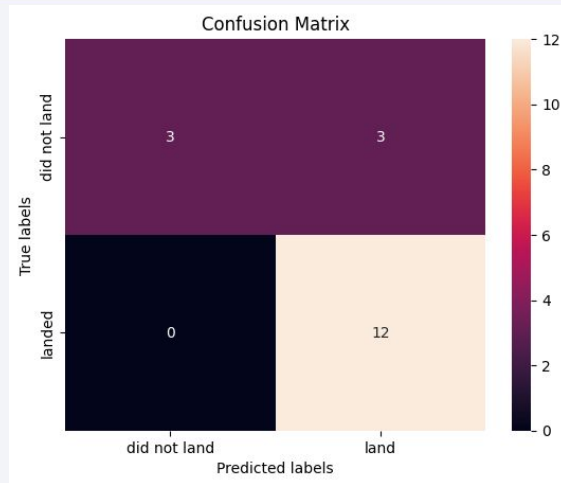
# Predictive Analysis (Classification)

# Classification Accuracy

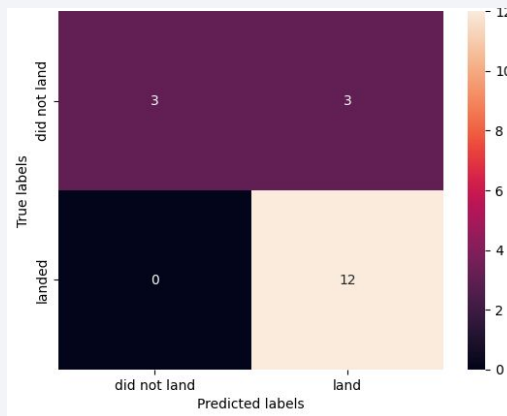
---

- Visualize the built model accuracy for all built classification models, in a bar chart
- Find which model has the highest classification accuracy

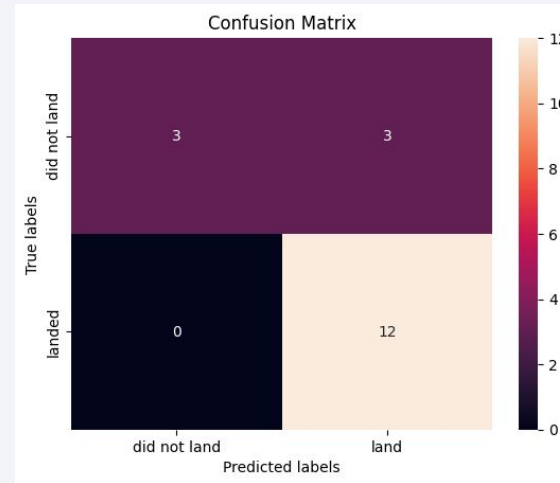
# Confusion Matrix



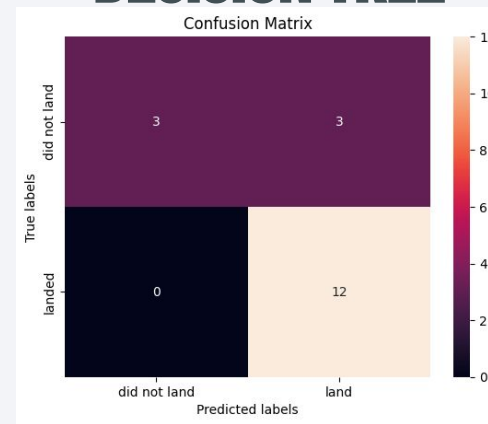
**LOGISTIC REGRESSION**



**SUPPORT VECTOR MACHINE**



**DECISION TREE**



**K NEAREST NEIGHBOUR**

ML Method	Accuracy Score (%)
Support Vector Machine	83.333333
Logistic Regression	83.333333
K Nearest Neighbour	83.333333
Decision Tree	83.333333

Examining the confusion matrix, we see that logistic regression can distinguish between the different classes. We see that the major problem is false positives.

# Conclusions

---

- Point 1
- Point 2
- Point 3
- Point 4
- ...

# Appendix

---

- Include any relevant assets like Python code snippets, SQL queries, charts, Notebook outputs, or data sets that you may have created during this project



Thank you!

