



# DATA SCIENCE

2025



# Outline

---

- Executive Summary
- Introduction
- Methodology
- Results
- Conclusion
- Appendix

# Executive Summary

---

- Summary of methodologies
- Python for Data Science
- Wrangling & Analysis
- Data Visualization (Matplotlib, Seaborn, Folium)
- Machine Learning Essentials
- SQL & Databases
- Capstone Project – End-to-end data science problem-solving.

# Introduction

---

- Project background and context The project aims to address data science
- , using data-driven methods to derive insights and solutions. Developed as part of the IBM Data Science Professional Certificate, this project is grounded in real-world data science
- ~~Objectives~~
- Analyze and visualize data to uncover hidden trends.
- Build predictive models to support business decision-making.
- Deliver actionable insights

Section 1

# Methodology

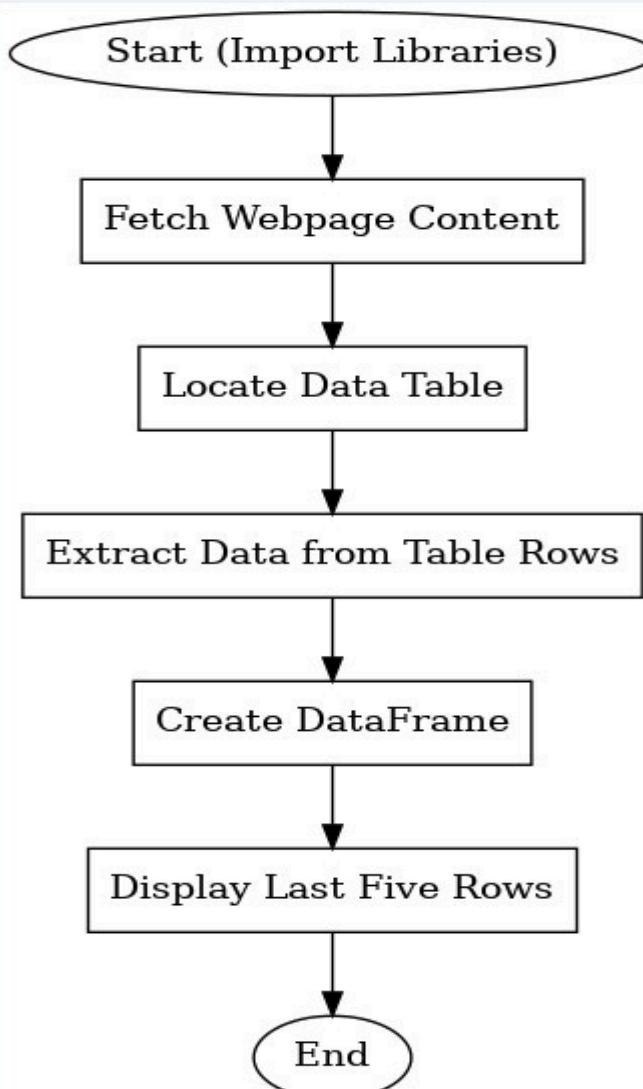
# Methodology

---

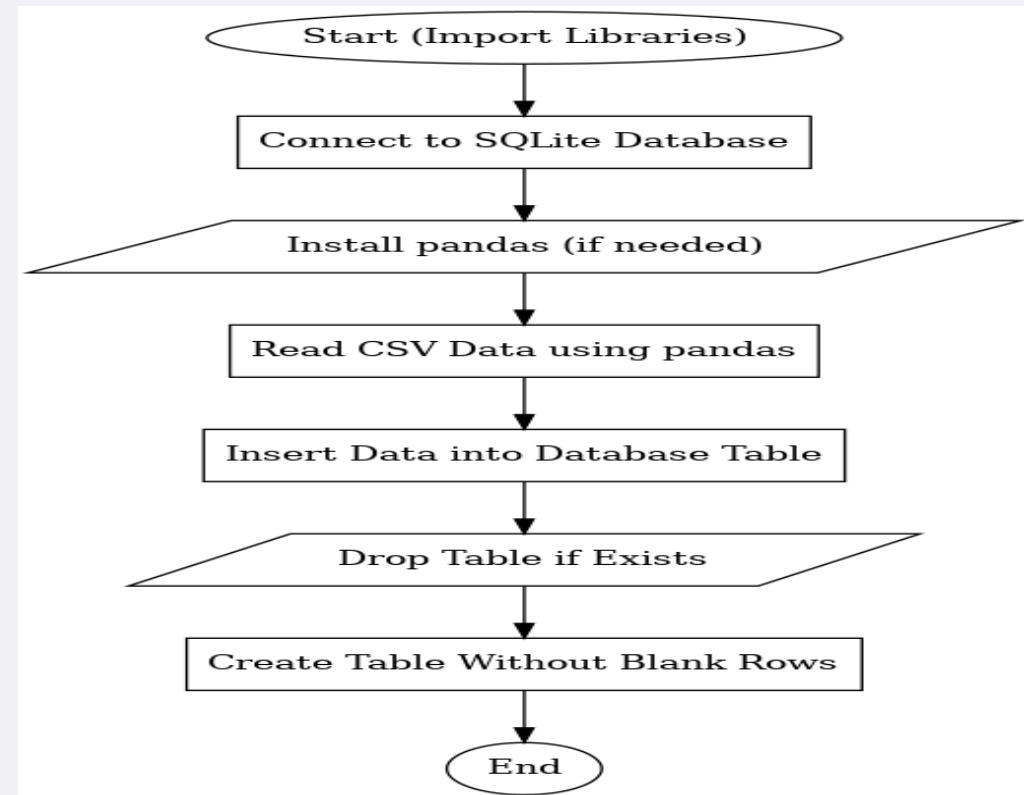
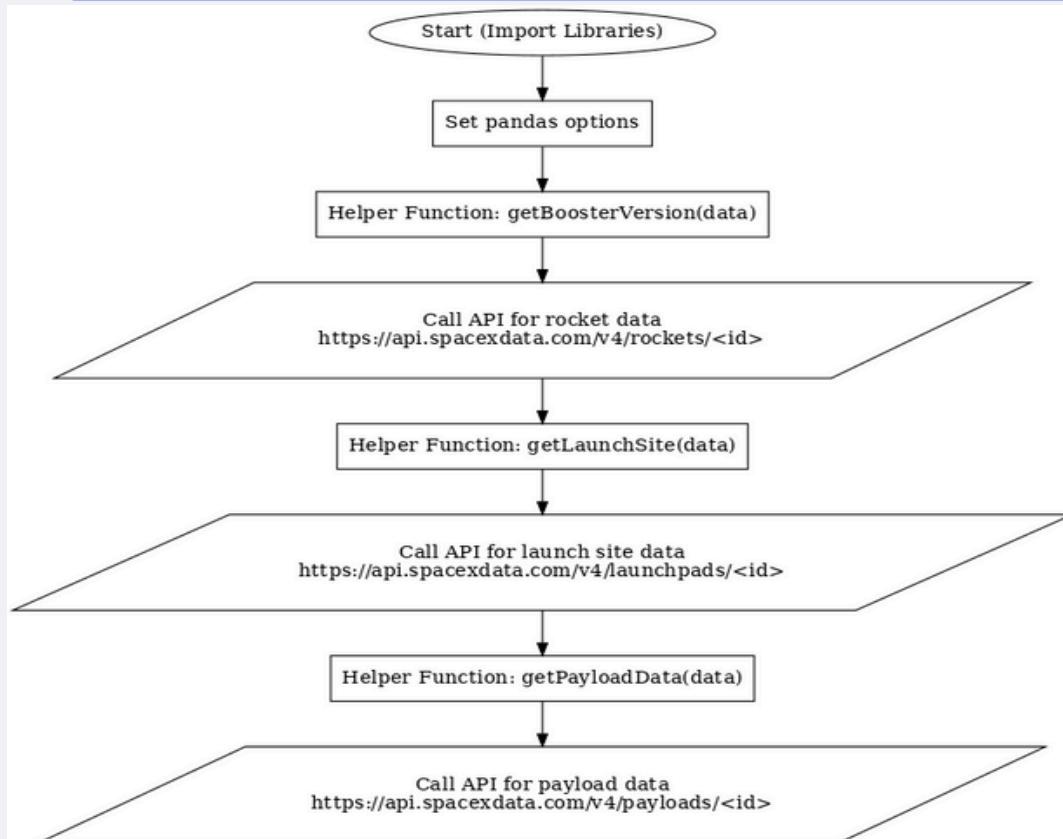
## Executive Summary

- Data collection methodology:
  - We use Wget function or Yfinance from yahoo for web scraping
- Perform datawrangling
  - We use pandas library to sort, arrange, clean and study data
  - Perform exploratory data analysis (EDA) using visualization andSQL
  - Perform interactive visual analytics using Folium andPlotlyDash
  - Perform predictive analysis using classification models
    - We use Scikitlearn we train test split data and use hyperparameter optimization with gridsearchCV

# Data Collection



# Data Collection – SpaceX API



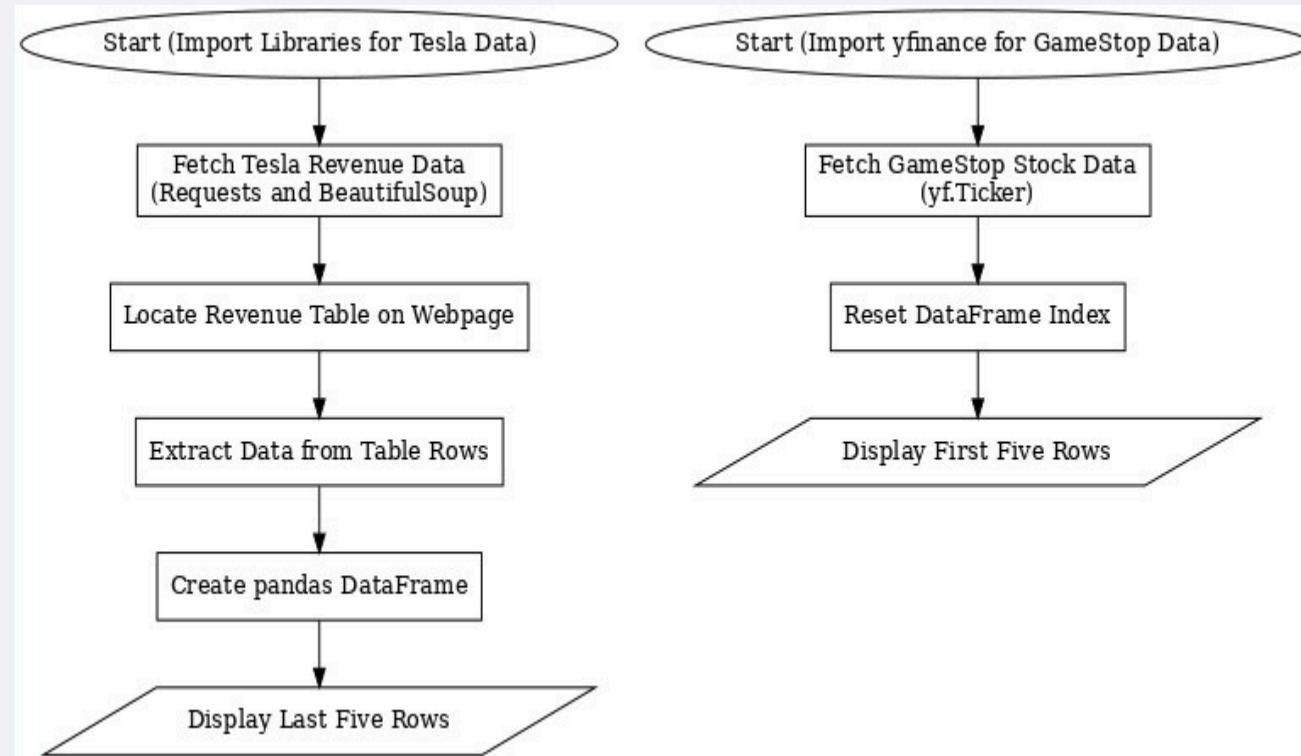
# Data Collection - Scraping

Start (Import Libraries)  
requests, BeautifulSoup, pandas.

FetchTesla RevenueData  
requests.get(), URL, headers, BeautifulSoup.

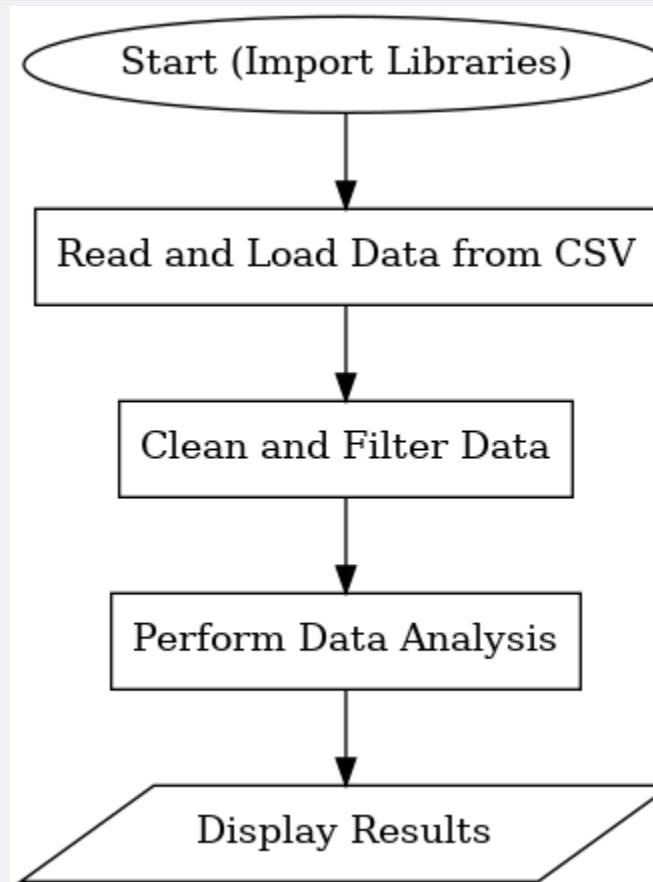
Locate RevenueTable:  
soup.find\_all(), table, historical\_data\_table.

ExtractData fromTable Rows: forloop,  
find\_all('td'), text.strip(), append.  
Createpandas DataFrame: pd.DataFrame(),  
data.



# Data Wrangling

---



# EDA with Data Visualization

Visualizing the Relationship between Flight Number and Launch Site Plot Type: Seaborn categorical plot (sns.catplot). Reason: This plot allows a clear comparison of how flights are distributed across various launch sites. Using hue="Class" distinguishes between successful and unsuccessful launches, providing additional insight into the distribution of outcomes by site and flight number.

Visualizing Success Rate by Orbit Type Plot Type: Bar plot (kind='bar'). Reason: Aggregating the success rate for each orbit type simplifies the comparison, allowing the viewer to quickly see which orbits have the highest or lowest success rates. A bar plot effectively communicates grouped summary statistics.

Visualizing Launch Success Yearly Trend Plot Type: Line plot (kind='line'). Reason: Line plots are ideal for time series data. This plot shows the trend of launch success rates over years, making it easy to spot improvements or regressions over time.

# EDA withSQL

---

- Create a Clean Table:

```
CREATE TABLE SPACEXTABLE AS SELECT * FROM SPACEXTBL WHERE Date IS NOT NULL
```

- Unique Launch Sites

```
SELECT DISTINCT Launch_Site FROM SPACEXTABLE
```

- Filtered Launch site Records

```
SELECT * FROM SPACEXTABLE WHERE Launch_Site LIKE 'CCA%' LIMIT 5
```

- Total Payload by NASA CRS Missions

```
SELECT SUM(PAYLOAD__MASS__KG_) as total_payload FROM SPACEXTABLE WHERE Customer LIKE '%NASA (CRS)%';
```

- Average Payload for Booster Version F9 v1.1

```
SELECT AVG(PAYLOAD__MASS__KG_) as avg_payload FROM SPACEXTABLE WHERE Booster_Version = 'F9 v1.1';
```

# Build an Interactive Map with Folium

---

- Circle Marker :
  - A circle was added to mark the location of NASA Johnson Space Center.
  - Features:
    - Radius: 1000 meters
    - Popuplabel: "NASA Johnson SpaceCenter"
- Text LabelMarker :
  - A Marker object was added with a DivIcon displaying the text "NASA JSC."
  - Purpose: Identify the Johnson Space Center with a visible label.
- Launch Outcome Markers :
  - Markers were added for each launch record in the dataset.
  - Features:
    - Color-coded based on success (green) or failure (red) of the launch.
    - Managed using the MarkerCluster plugin to reduce map clutter for overlapping points.
- Mouse Position Plugin :
  - Added to the map to display coordinates of points under the mouse cursor.
  - Purpose: To help identify locations (e.g., railways, highways, coastline) near the launch sites for further proximity analysis.
- Distance Markers :
  - Markers were created to indicate distances from launch sites to points of interest (e.g., closest coastline, city, or highway).
  - Features:
    - Custom HTML DivIcon used to display the distance.

# Build a Dashboard with Plotly Dash

---

## 1. Dropdown Menu (Launch Site Selection)

- o Purpose: Allows the user to select a specific launch site or view data for all sites.
  - o Interaction: Triggers updates in both the pie chart and scatter plot.
- 

## 2. Pie Chart (Launch Success Overview)

- o Purpose: Displays total successful launches for all sites or success vs. failure for a selected site.
  - o Data:
    - When "All Sites" is selected, it shows success counts per site.
    - When a specific site is selected, it shows a breakdown of success vs. failure launches.
  - o Dynamic Interaction: Updates based on the selected launch site from the dropdown.
- 

## 3. Payload Range Slider

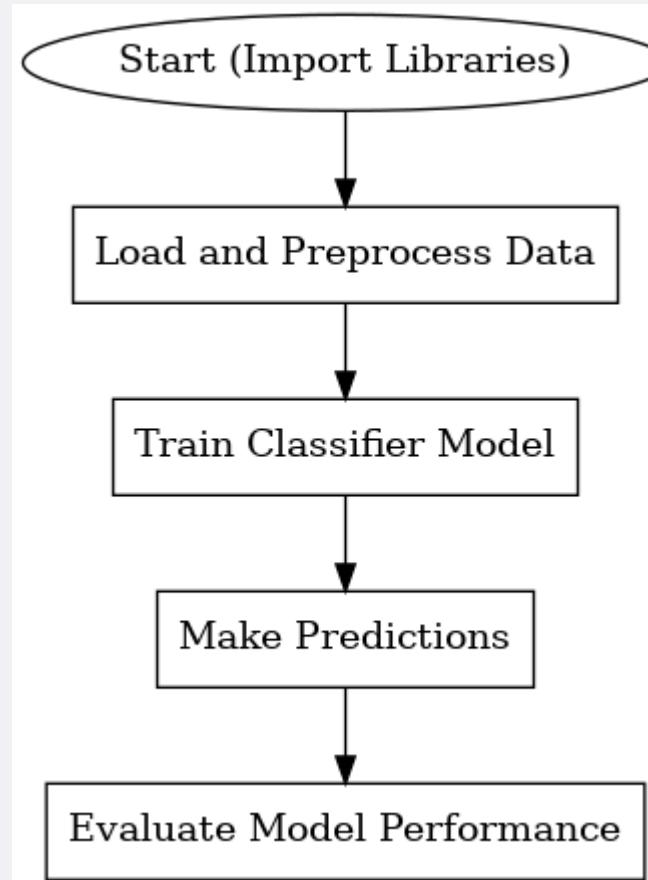
- o Purpose: Enables users to filter data based on the payload mass (in kilograms).
  - o Interaction: The slider adjusts the payload range for the scatter plot.
- 

## 4. Scatter Plot (Payload vs. Launch Success)

- o Purpose: Visualizes the correlation between payload mass and launch success rate.
  - o Data:
    - X-axis: Payload Mass (kg)
    - Y-axis: Launch outcome (0 = Failure, 1 = Success)
    - Color: Booster version category
  - o Dynamic Interaction: Updates based on both the selected launch site and payload range
-

# Predictive Analysis (Classification)

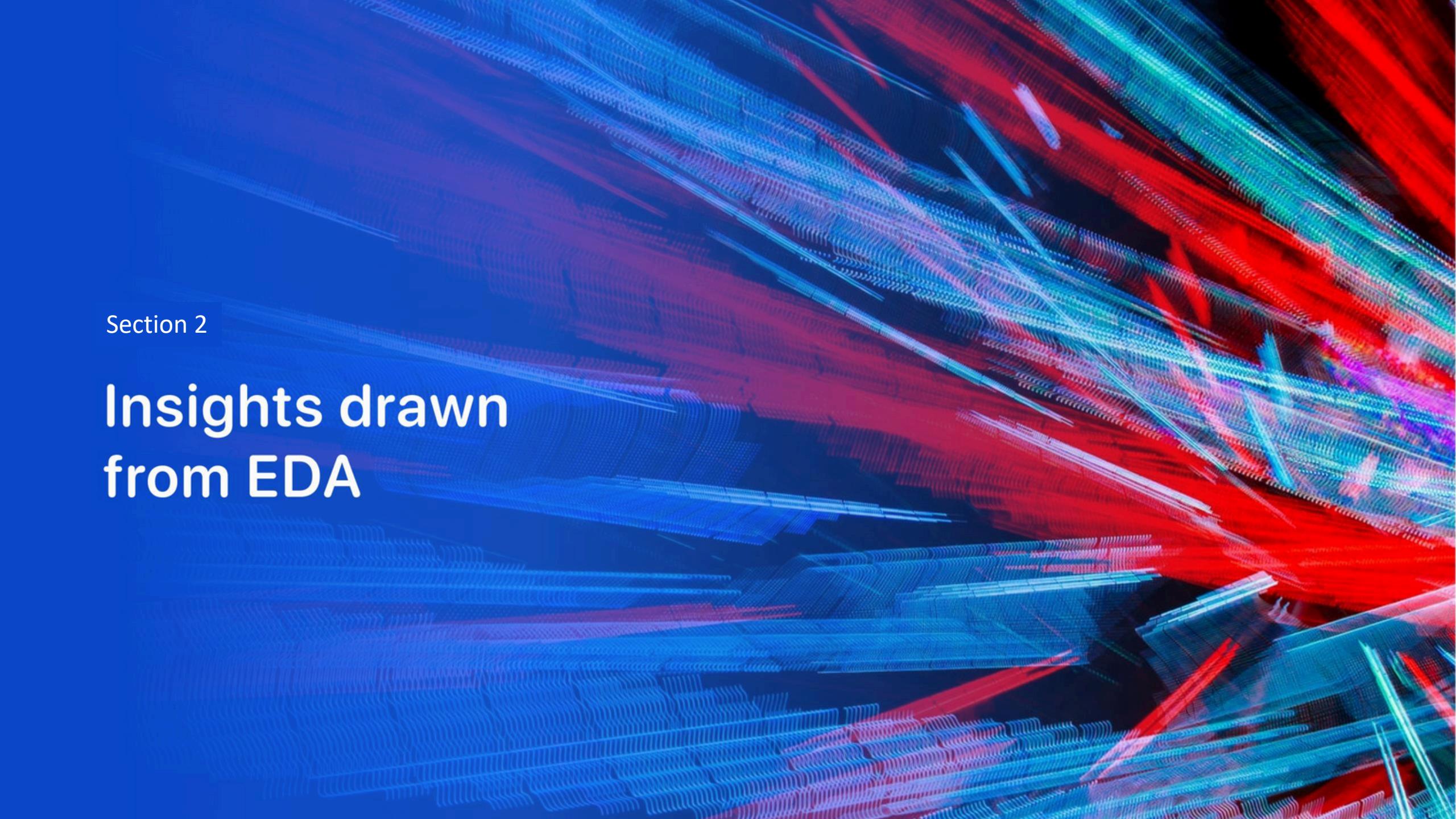
---



# Results

---

- Data Loading and Preparation: Data was read and cleaned to remove missing values or inconsistencies.
- Model Training:
- Classifiers were trained on the processed dataset to predict weather patterns.
- Evaluation Metrics: The model's accuracy, precision, and recall were calculated to assess performance.
- Insights and Visualization: Charts or tables likely summarized key weather-related insights, such as temperature trends or humidity effects on outcomes.

The background of the slide features a dynamic, abstract pattern of glowing particles. These particles are arranged in numerous wavy, flowing lines that create a sense of motion. The colors used are primarily shades of blue, red, and green, which are bright and stand out against the dark, almost black, background. The overall effect is reminiscent of a digital or quantum simulation visualization.

Section 2

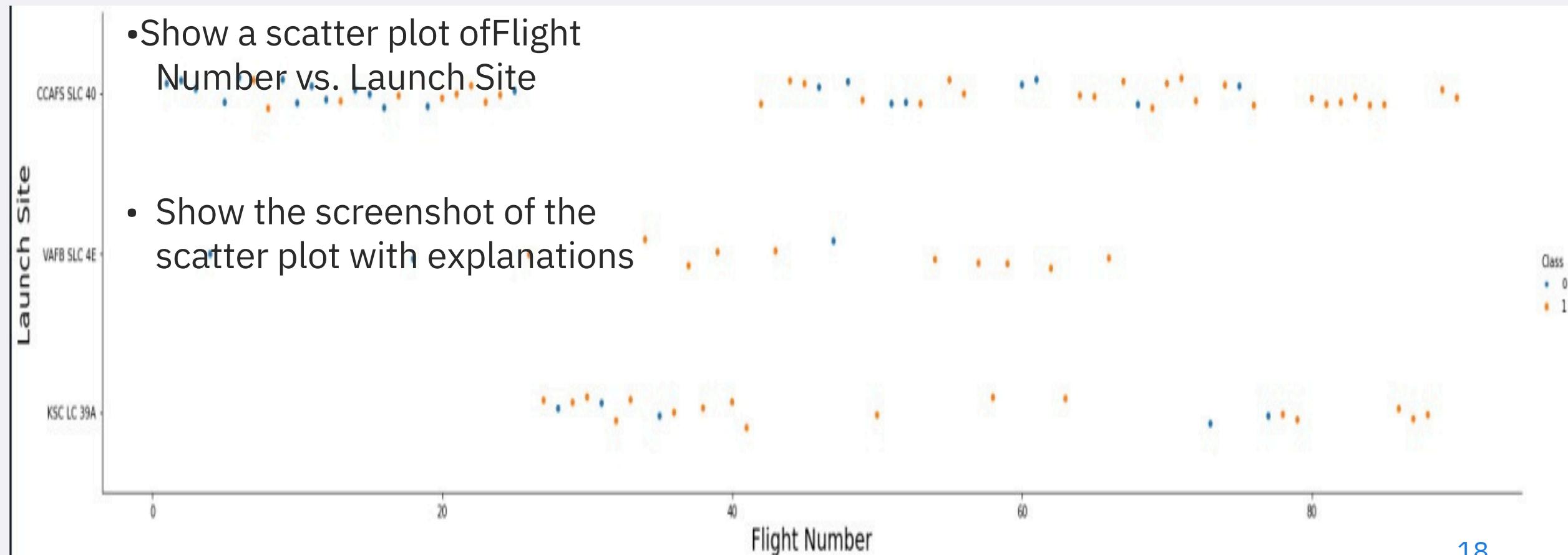
## Insights drawn from EDA

# Flight Number vs. LaunchSite

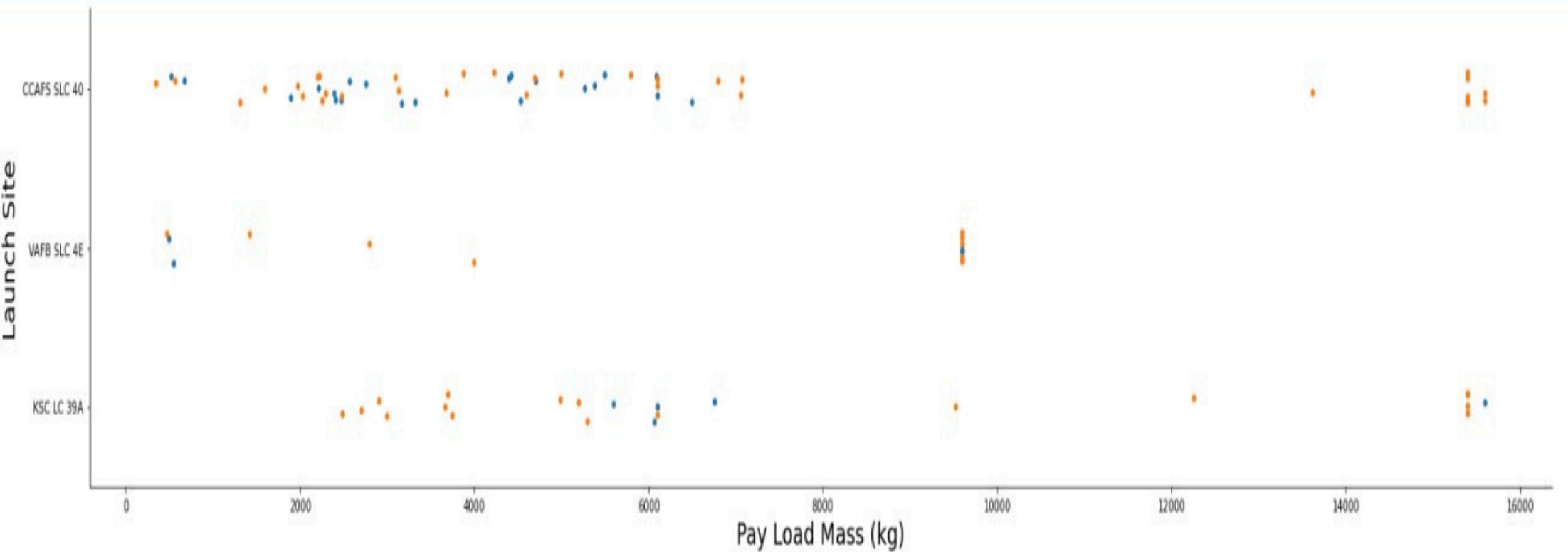
---

- Show a scatter plot of Flight Number vs. Launch Site

- Show the screenshot of the scatter plot with explanations

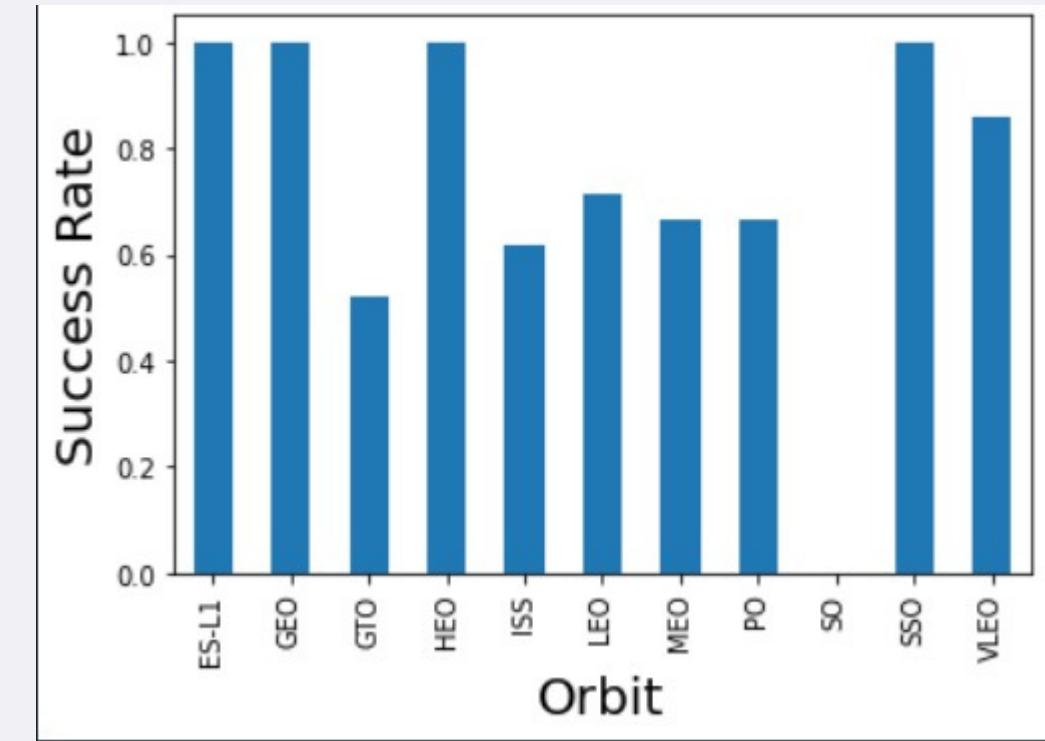
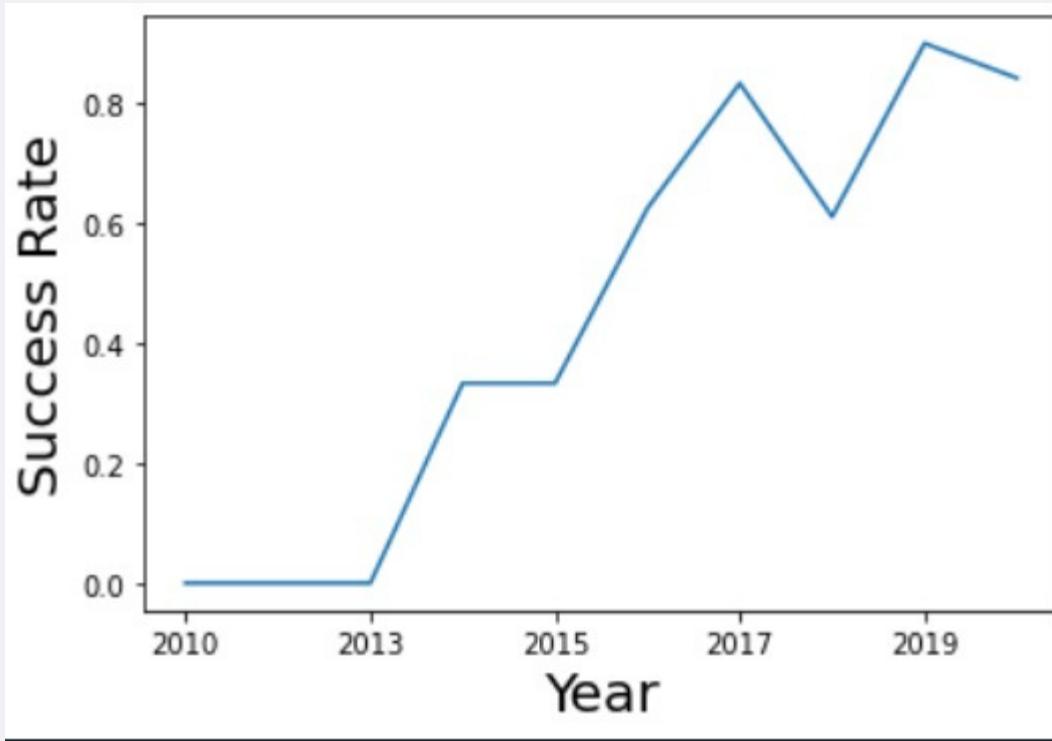


# Payload vs. Launch Site



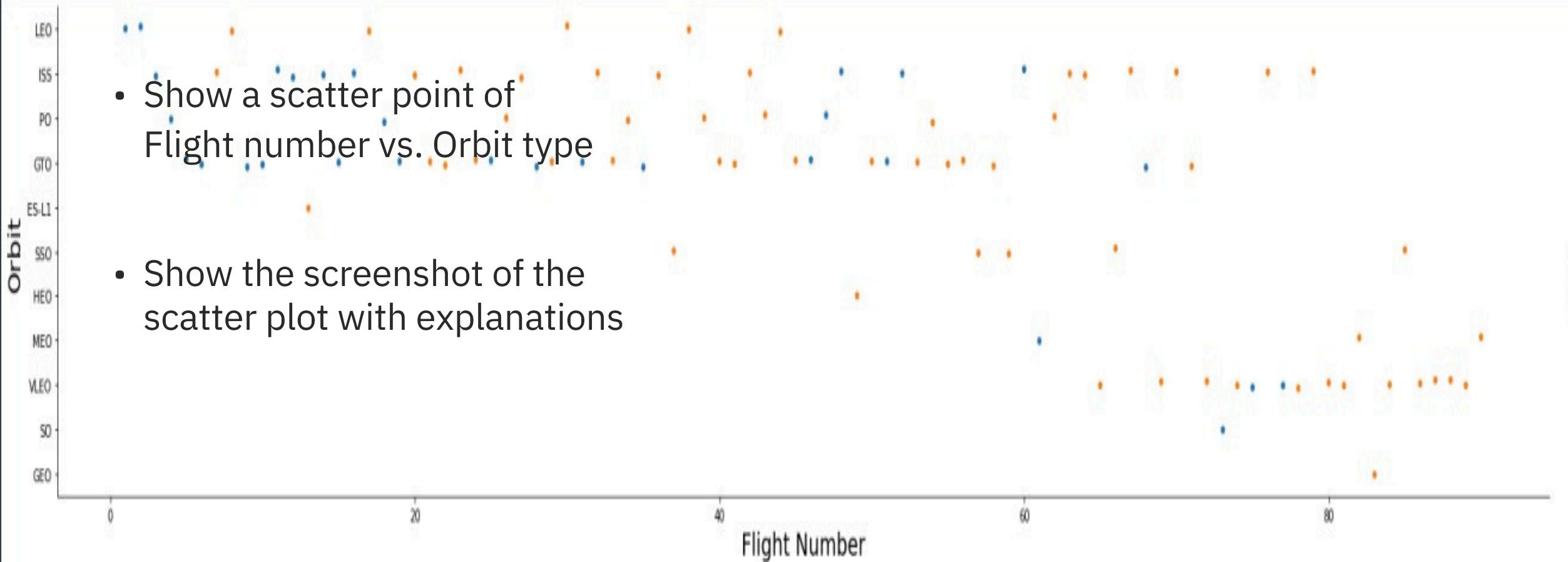
# Success Rate vs. Orbit Type

---



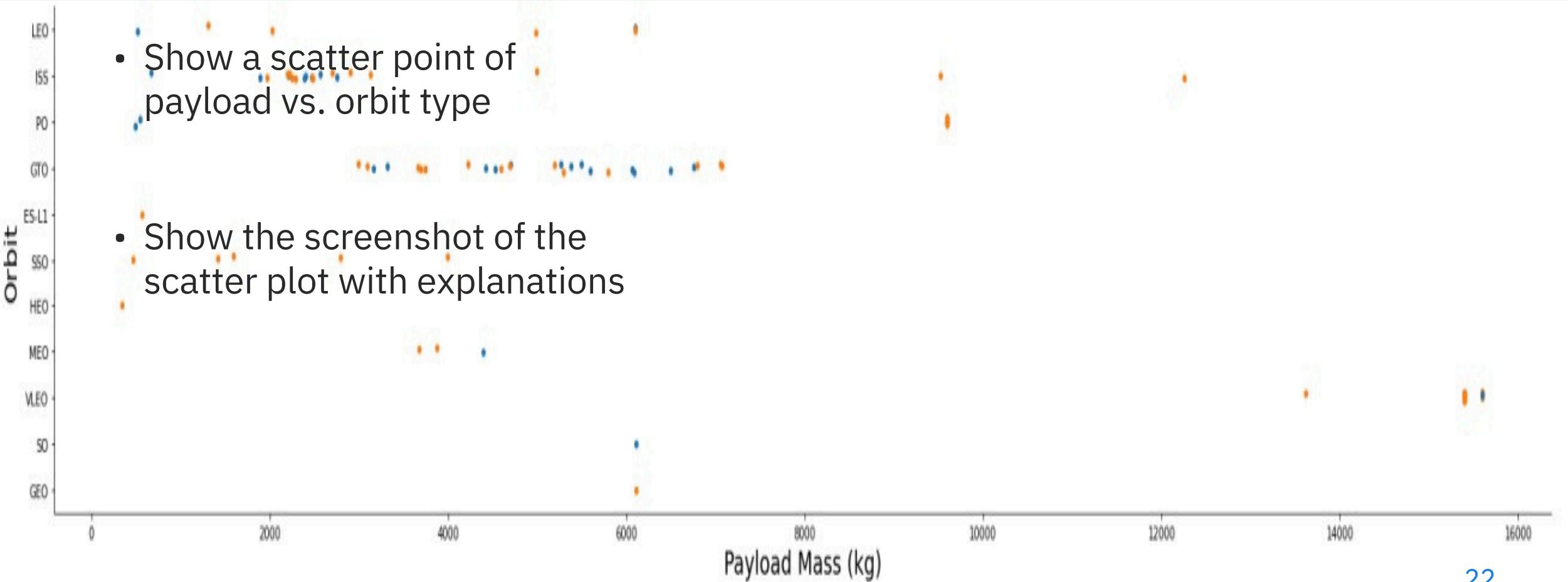
- There was learning about orbits
- Show the screenshot of the scatter plot with explanations

# Flight Number vs. Orbit Type



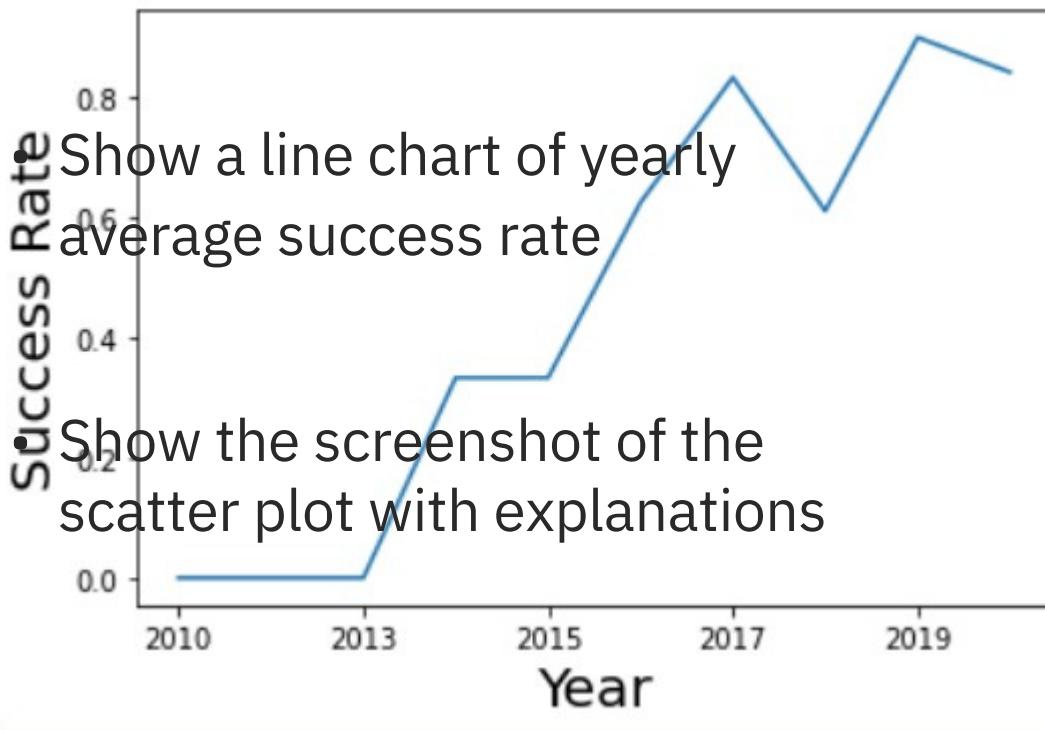
# Payload vs. OrbitType

---



# Launch Success Yearly Trend

---



# All Launch Site Names

---

- Find the names of the unique launch sites
- SELECT DISTINCT Launch\_Site FROM SPACEXTABLE;
- Launch Site
- CCAFS LC-401
- VAFB SLC-4E2
- KSC LC-39A3
- CCAFS SLC-40

# Launch Site Names Begin with 'CCA'

---

- Find 5 records where launch sites begin with `CCA`
- ```
query_cca_records = "SELECT * FROM SPACEXTABLE WHERE Launch_Site
LIKE 'CCA%' LIMIT 5"
cursor.execute(query_cca_records)
cca_records =
cursor.fetchall()
```

|   | Date       | Time (UTC) | Booster_Version | Launch_Site | Payload                              | PAYOUTLOAD_MASS_KG | Orbit     | Customer        | Mission_Outcome | Landing_Outcome     |
|---|------------|------------|-----------------|-------------|--------------------------------------|--------------------|-----------|-----------------|-----------------|---------------------|
| 1 | 2010-06-04 | 18:45:00   | F9 v1.0 B0003   | CCAFS LC-40 | Dragon Spacecraft Qualification Unit | 0                  | LEO       | SpaceX          | Success         | Failure (parachute) |
| 2 | 2010-12-08 | 15:43:00   | F9 v1.0 B0004   | CCAFS LC-40 | Dragon demo flight C1, two CubeSa... | 0                  | LEO (ISS) | NASA (COTS) NRO | Success         | Failure (parachute) |
| 3 | 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          |
| 4 | 2012-10-08 | 0:35:00    | F9 v1.0 B0006   | CCAFS LC-40 | SpaceX CRS-1                         | 500                | LEO (ISS) | NASA (CRS)      | Success         | No attempt          |
| 5 | 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

---

- # Calculate the total payload mass carried by all boosters
- `total_payload_mass = df['PAYLOAD_MASS__KG_'].sum()`
- 619967 Kg

# Average Payload Mass by F9 v1.1

---

- Calculate the average payload mass carried by booster version F9 v1.1
- `average_payload_f9_v1_1 = df_new[df_new['Booster_Version'] == 'F9 v1.1']['PAYLOAD__MASS__KG_'].mean()`
- `average_payload_f9_v1_1 = 2928.4 kg`

# First Successful Ground Landing Date

---

- Find the dates of the first successful landing outcome on ground pad
- `first_success_ground_pad_date = df_new[df_new['Landing_Outcome'] == 'Success (ground pad)']['Date'].min()`
- `first_success_ground_pad_date = '2015-12-22'`

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

---

- List the names of boosters which have successfully landed on drone ship and had payload mass greater than 4000 but less than 6000
- ```
boosters_success_drone_ship = df_new[(df_new['Landing_Outcome'] == 'Success (drone ship)') & (df_new['PAYLOAD_MASS__KG_'] > 4000) & (df_new['PAYLOAD_MASS__KG_'] < 6000)]['Booster_Version'].unique()
```
- The boosters that have successfully landed on a drone ship and carried a payload mass between 4000 and 6000 kg are:
  - 1.F9 FT B1022**
  - 2.F9 FT B1026**
  - 3.F9 FT B1021.2**
  - 4.F9 FT B1031.2**

## Total Number of Successful and Failure Mission Outcomes

---

- Calculate the total number of successful and failure mission outcomes
- `mission_outcomes_counts = df_new['Mission_Outcome'].value_counts()`
  
- `mission_outcomes_counts`
- Success 98
- Failure (in flight) 1
- Success (payload status unclear) 1
- Success 1

# Boosters Carried Maximum Payload

---

- List the names of the booster which have carried the maximum payload mass

```
max_payload_mass = df_new['PAYLOAD__MASS__KG_'].max()
```

```
# Filter the data to find boosters that carried the maximum  
payload mass
```

- ```
boosters_with_max_payload =  
df_new[df_new['PAYLOAD__MASS__KG_'] ==  
max_payload_mass]['Booster_Version'].unique()
```

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

# 2015 Launch Records

---

- List the failed landing\_outcomes in drone ship, their booster versions, and launch site names for in year 2015
- # Filter the data for failed drone ship landings in the year 2015
- failed\_drone\_ship\_2015 = df\_new[
  - (df\_new['Landing\_Outcome'].str.contains('Failure', na=False)) &
  - (df\_new['Landing\_Outcome'].str.contains('ship', na=False)) &
  - (df\_new['Date'].str.startswith('2015'))
- ][['Landing\_Outcome', 'Booster\_Version', 'Launch\_Site']]
- 13 Failure (drone ship)
- F9 v1.1 B1012 CCAFS LC-40
- F9 v1.1 B1015 CCAFS LC-40

## 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
- `ranked_landing_outcomes = df_new[  
 (df_new['Date'] >= '2010-06-04') & (df_new['Date'] <= '2017-03-20')  
 ]['Landing_Outcome'].value_counts().reset_index()  
 # Rename columns for clarity  
 ranked_landing_outcomes.columns = ['Landing_Outcome', 'Count']  
 # Display the results sorted in descending order  
 ranked_landing_outcomes`

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 the United States appears. There are also larger clusters of lights in South America and Europe. The atmosphere of the Earth is visible as a thin blue layer, and the horizon line is clearly defined.

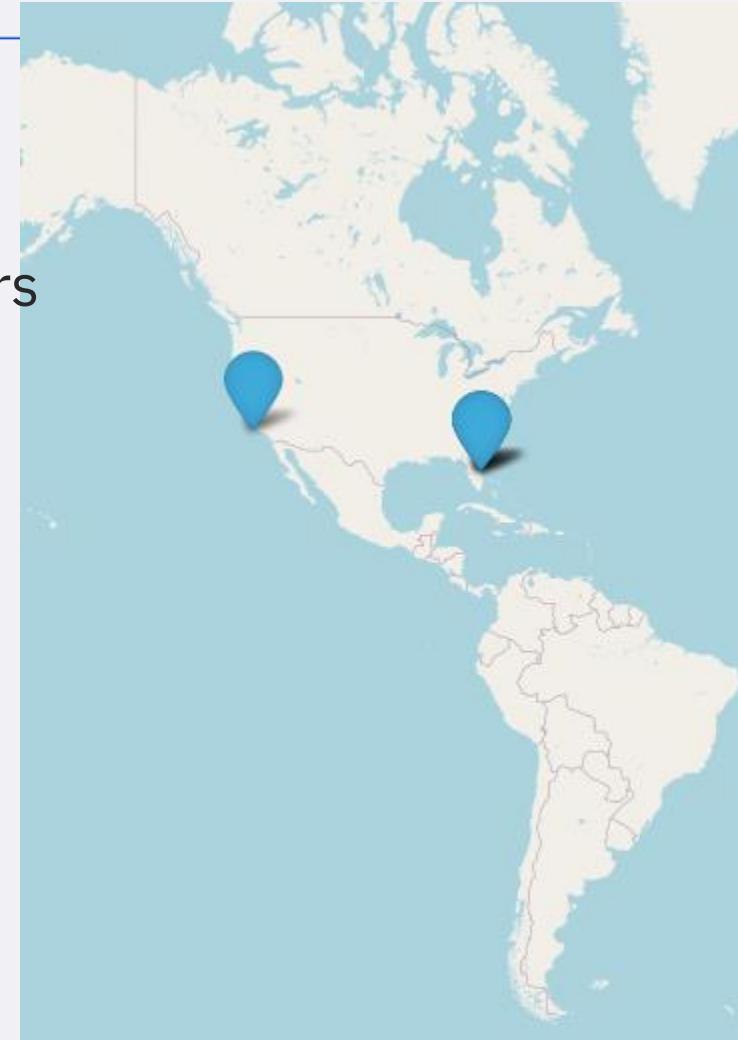
Section 3

# Launch Sites Proximities Analysis

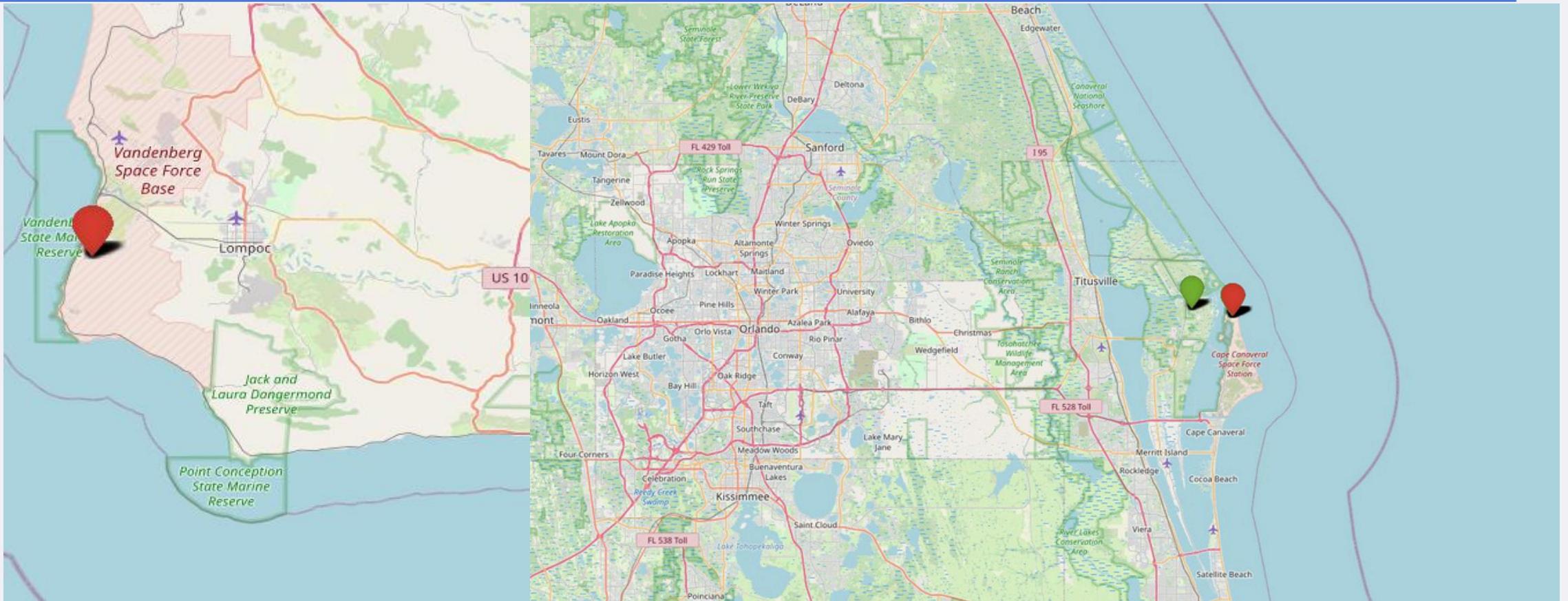
# Launch sites map

---

- Explore the generated folium map and make a proper screenshot to include all launch sites' location markers on a global map
- Each marker on the map represents a SpaceX launch site. The icon used is a blue rocket, symbolizing a launch location.

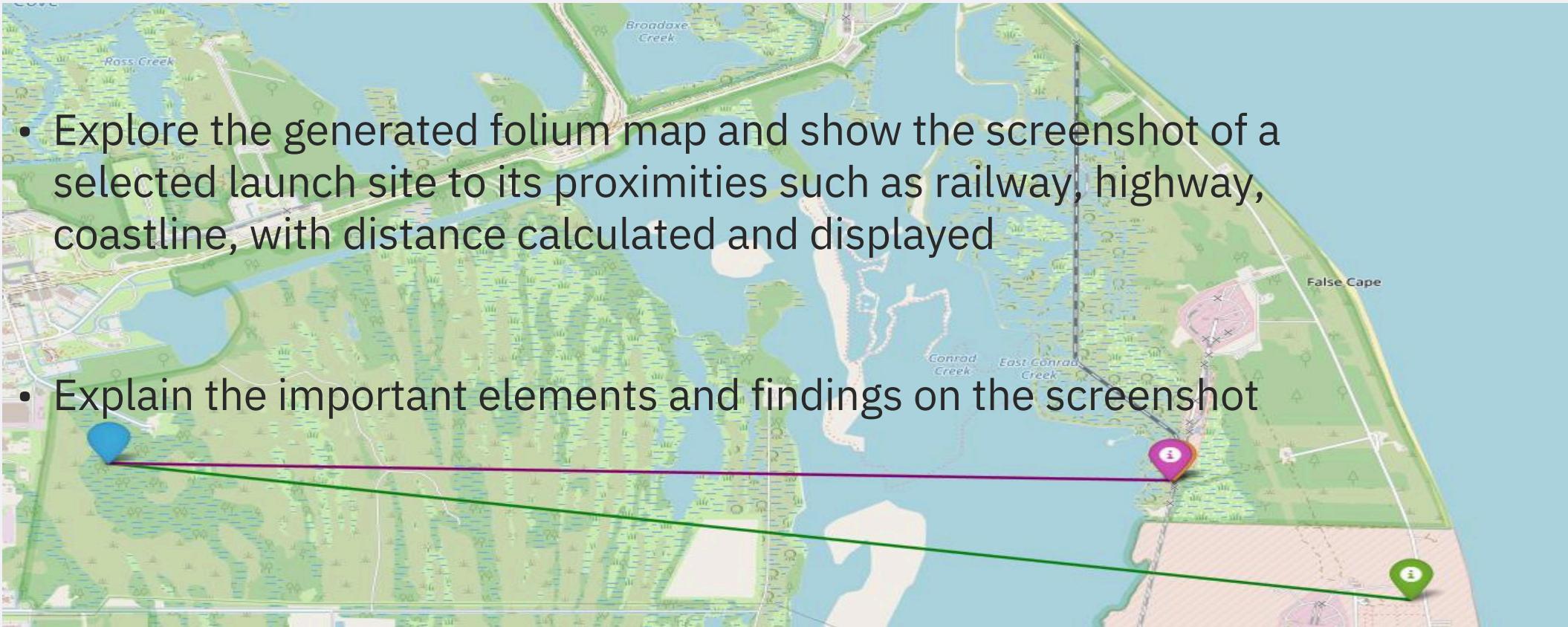


# Launch Sites 2

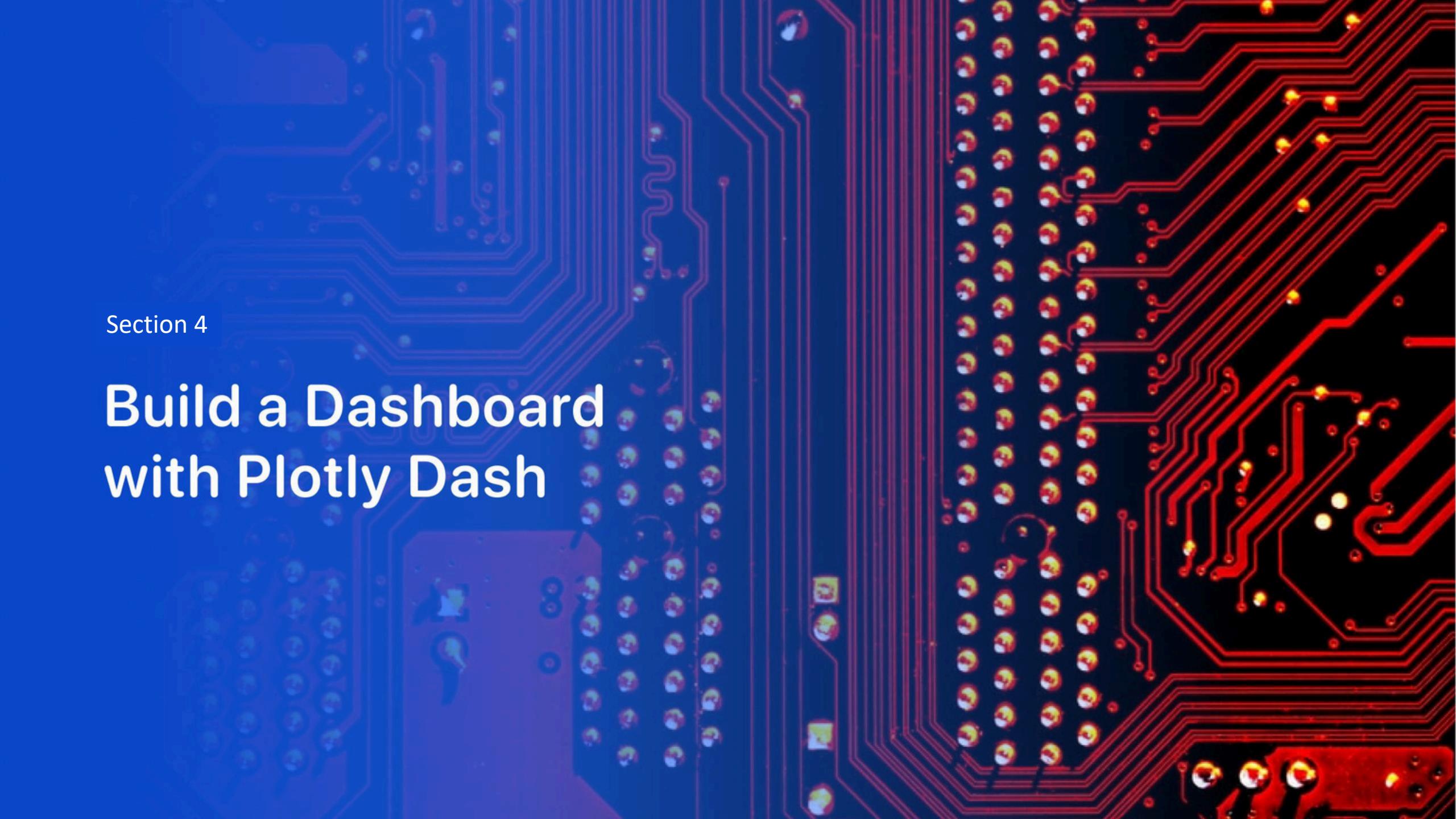


Each launch site marker is color-coded to indicate the outcome of the launch:  
Green Marker: Represents a successful launch.  
Red Marker: Represents a failed launch.

# Launch sites Risks



The blue marker with the rocket icon indicates the exact location of the Kennedy Space Center (KSC LC-39A) launch site. Clicking on the marker shows a popup with the name of the launch site. Proximity Markers: Green Marker (Coastline): Represents the closest point on the coastline. Distance from the launch site is displayed in kilometers. Orange Marker (Highway): Represents the nearest highway to the launch site. Distance from the launch site is shown on the marker



Section 4

# Build a Dashboard with Plotly Dash

# Space X dashboard

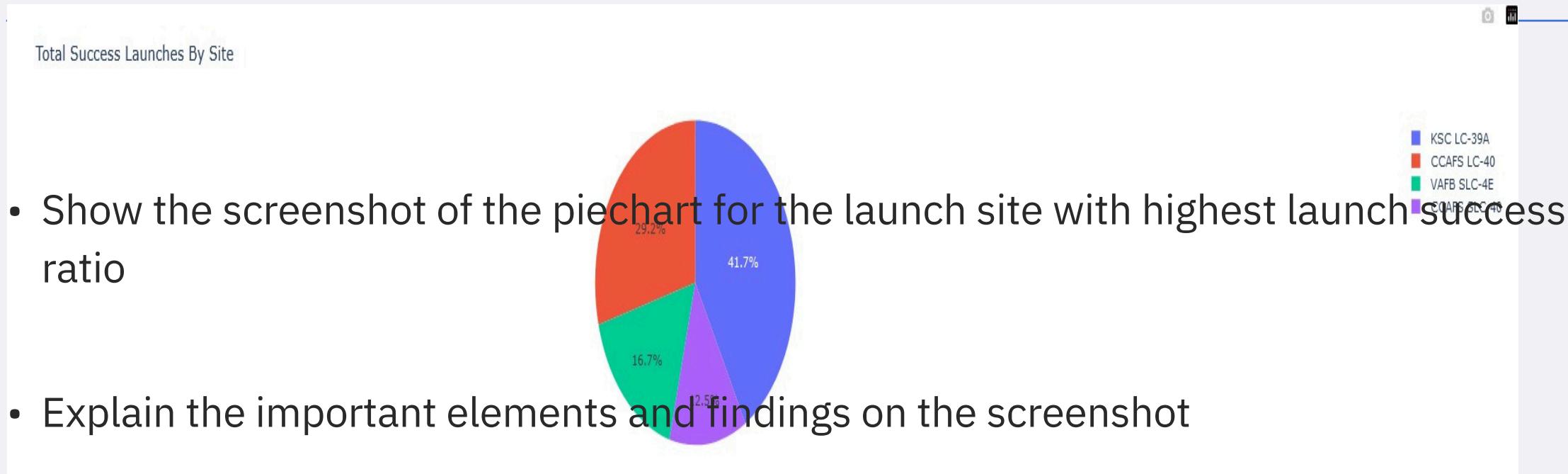
Show the screenshot of launch success count for all sites, in a piechart

- Explain the important elements and findings on the screenshot



KSC LC-39 A mostsuccessful

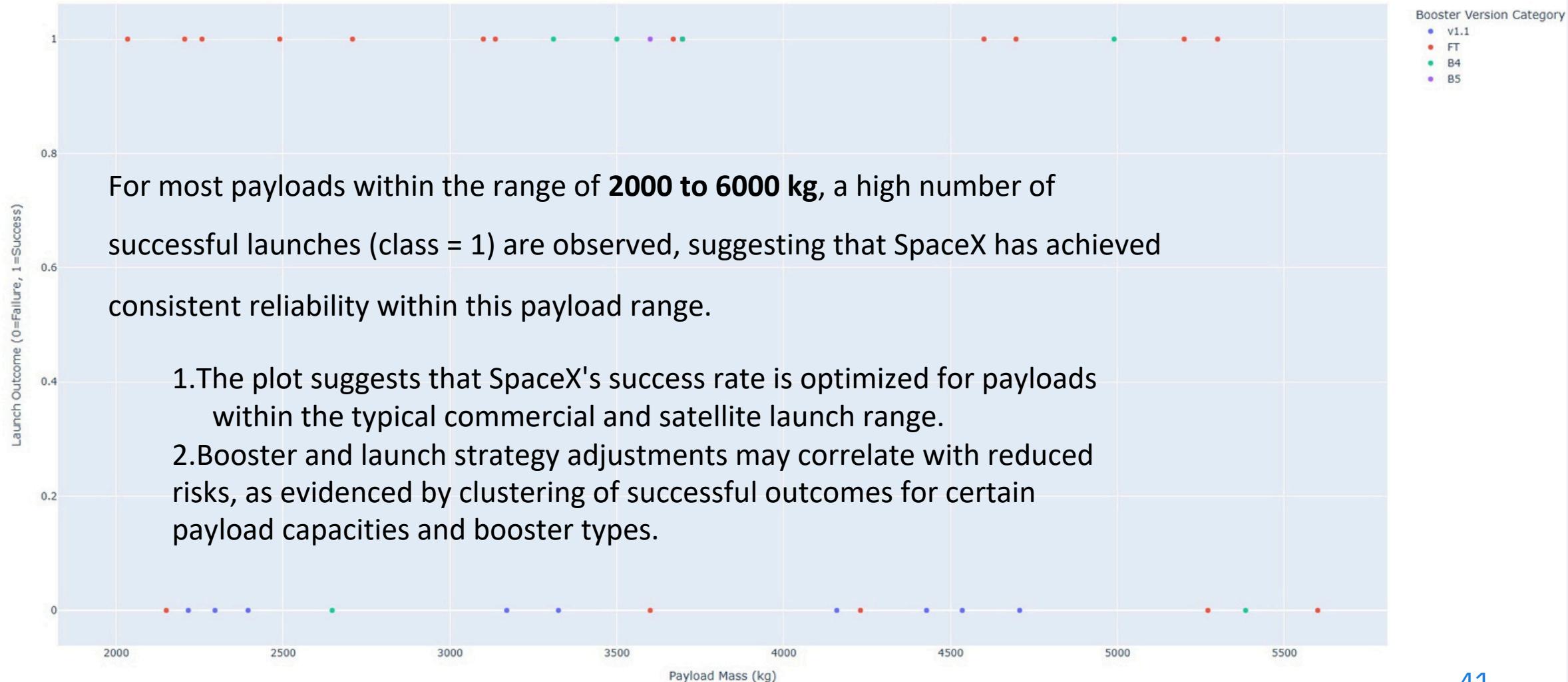
# Space X Dashboard 2



KSC LC-39 A mostsuccessful

# Payload vs Success all sites

Payload vs. Success for All Sites (Filtered by Payload Range)



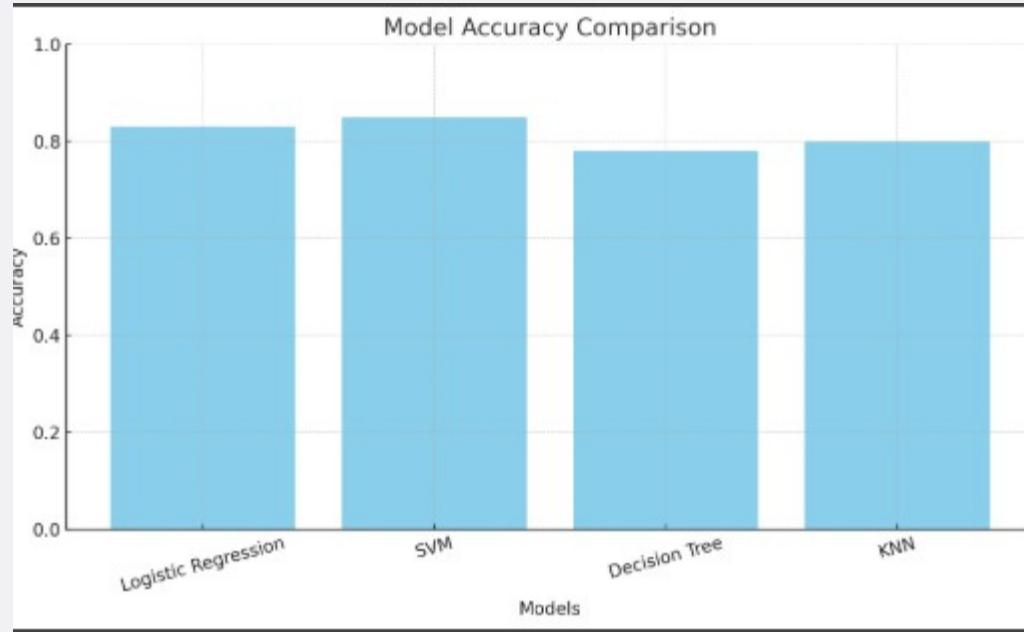
The background of the slide features a dynamic, abstract design. It consists of several thick, curved lines in shades of blue and yellow, creating a sense of motion and depth. The lines curve from the top right towards the bottom left, with some lines being more prominent than others. The overall effect is reminiscent of a tunnel or a high-speed train track.

Section 5

# Predictive Analysis (Classification)

# Classification Accuracy

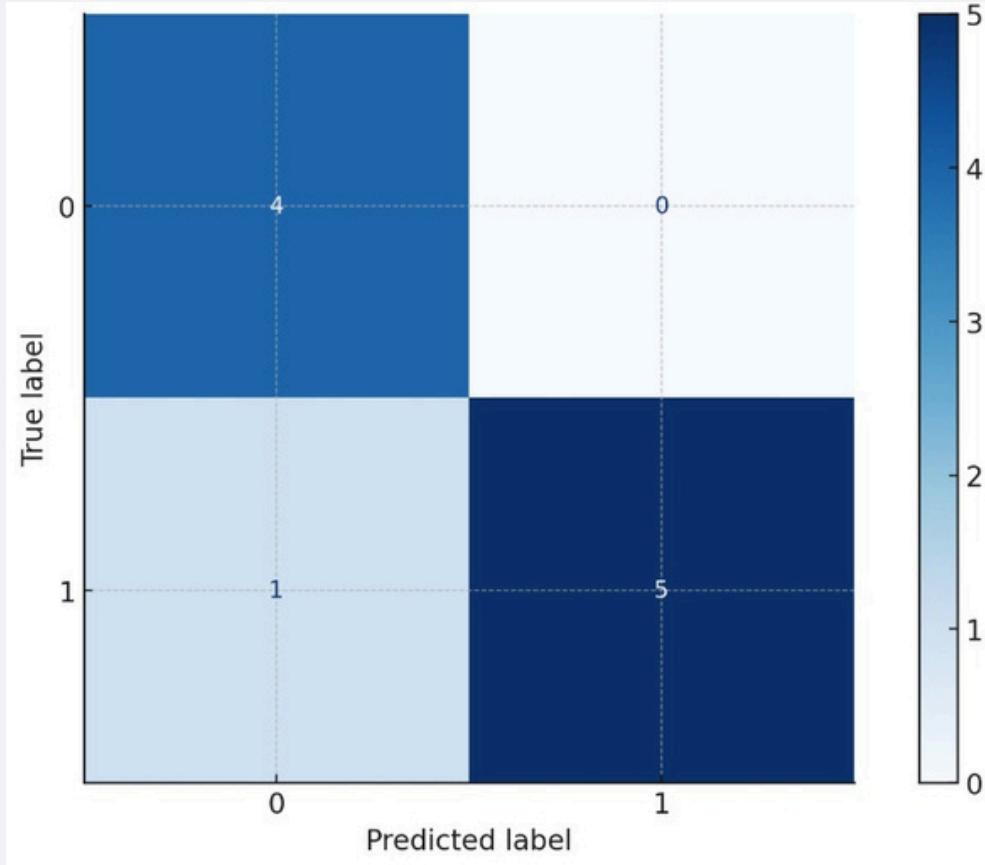
---



The model with the highest accuracy is **Support Vector Machine (SVM)** with an accuracy of **0.85**.

# Confusion Matrix

---



True Negative (TN): 4  
Correct predictions where the actual outcome was 0 (failure) and the model predicted 0.

False Positive (FP): 0  
Incorrect predictions where the actual outcome was 0 but the model predicted 1 (success).

False Negative (FN): 1  
Incorrect predictions where the actual outcome was 1 (success) but the model predicted 0.

True Positive (TP): 5  
Correct predictions where the actual outcome was 1 (success) and the model predicted 1.

# 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!

