

Winning Space Race with Data Science

Millicent Goodwin
November 30, 2024



Outline

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

Executive Summary

- Comprehensive analysis of SpaceX Falcon 9 launch data to predict first stage landing success
- Data collection from SpaceX API and web scraping of Wikipedia
- Exploratory Data Analysis using Python, SQL, and visualization tools
- Development of machine learning models achieving 83% accuracy in predicting landing outcomes
- Geospatial analysis revealing launch site characteristics and success patterns

Introduction

Background

- SpaceX revolutionized space industry with reusable Falcon 9 first stage
- Cost reduction from \$62M to ~\$30M per launch due to reusability
- Crucial to predict landing success for operational planning

Project Goals

- Predict first stage landing success using machine learning
- Analyze factors influencing landing outcomes
- Develop interactive visualizations for stakeholder insights

Section 1

Methodology

Methodology

Data collection methodology:

- SpaceX API: Used REST API calls to gather detailed information about past rocket launches including payload mass, booster versions, landing success, and more
- Web Scraping: Extracted historical Falcon 9 launch records from Wikipedia to validate and supplement API data.

Perform data wrangling:

- Replaced missing values through imputation, especially in landing outcomes and payload masses
- Processed datetime features for launch dates
- Converted categorical variables into numerical formats for analysis

Methodology

Exploratory Data Analysis (EDA):

- Performed visualization analysis using Python libraries
- Executed SQL queries to gain insights into launch patterns
- Created geospatial visualizations with Folium for launch site analysis
- Built interactive dashboards using Plotly Dash

Machine Learning Implementation:

- Developed classification models including:
 - Logistic Regression
 - K-Nearest Neighbors (KNN)
 - Random Forest
 - Support Vector Machines (SVM)

Methodology

- Evaluated models using metrics like accuracy, precision, and recall
- Selected the best performing model for landing success prediction

This comprehensive methodology enabled the development of accurate predictive models for SpaceX Falcon 9 first stage landing success, providing valuable insights for cost estimation and operational planning.

Data Collection

The dataset used in this project was gathered in two main ways:

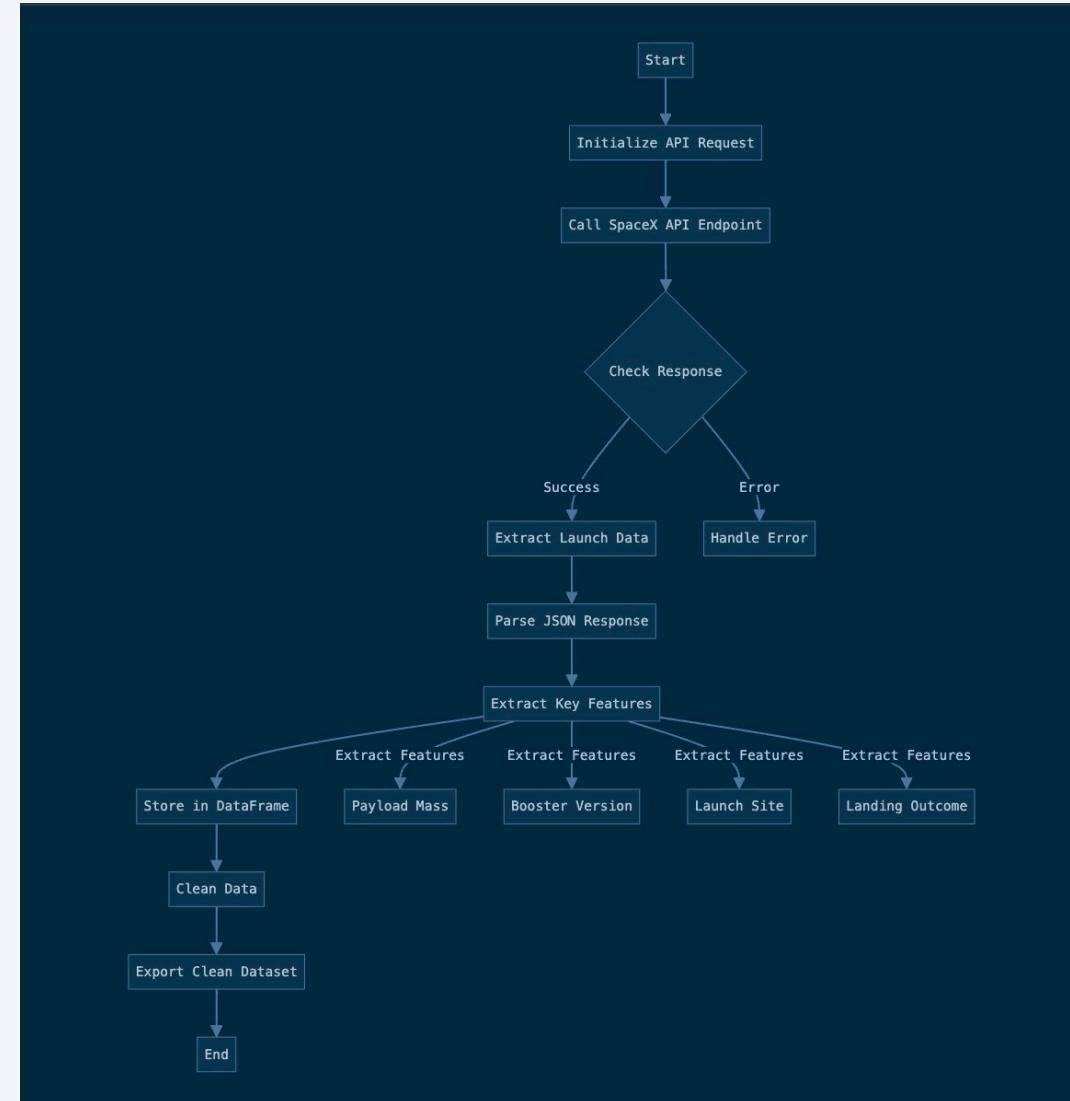
- 1. SpaceX API:** The SpaceX API provides detailed information about past rocket launches. The data includes crucial features such as payload mass, booster versions, landing success, and more. API calls were made to retrieve launch-related data, and auxiliary functions were defined to collect additional information about rocket boosters, launch sites, payloads, and landing cores.
- 2. Web Scraping:** The second was gathered by scraping historical Falcon 9 launch records from the [Wikipedia](#) page. Using BeautifulSoup, relevant columns such as flight numbers, payload masses, and landing outcomes were extracted into a clean dataset.

Data Collection – SpaceX API

GitHub Repository Reference:
<https://github.com/phdj91/Applied-Data-Science-Capstone>

Location of SpaceX API calls notebook:

- Repository: Applied-Data-Science-Capstone
- Notebook path: Applied-Data-Science-Capstone/spacex_api_data_collection.ipynb
- Contains complete code cells and outcomes for API data collection process



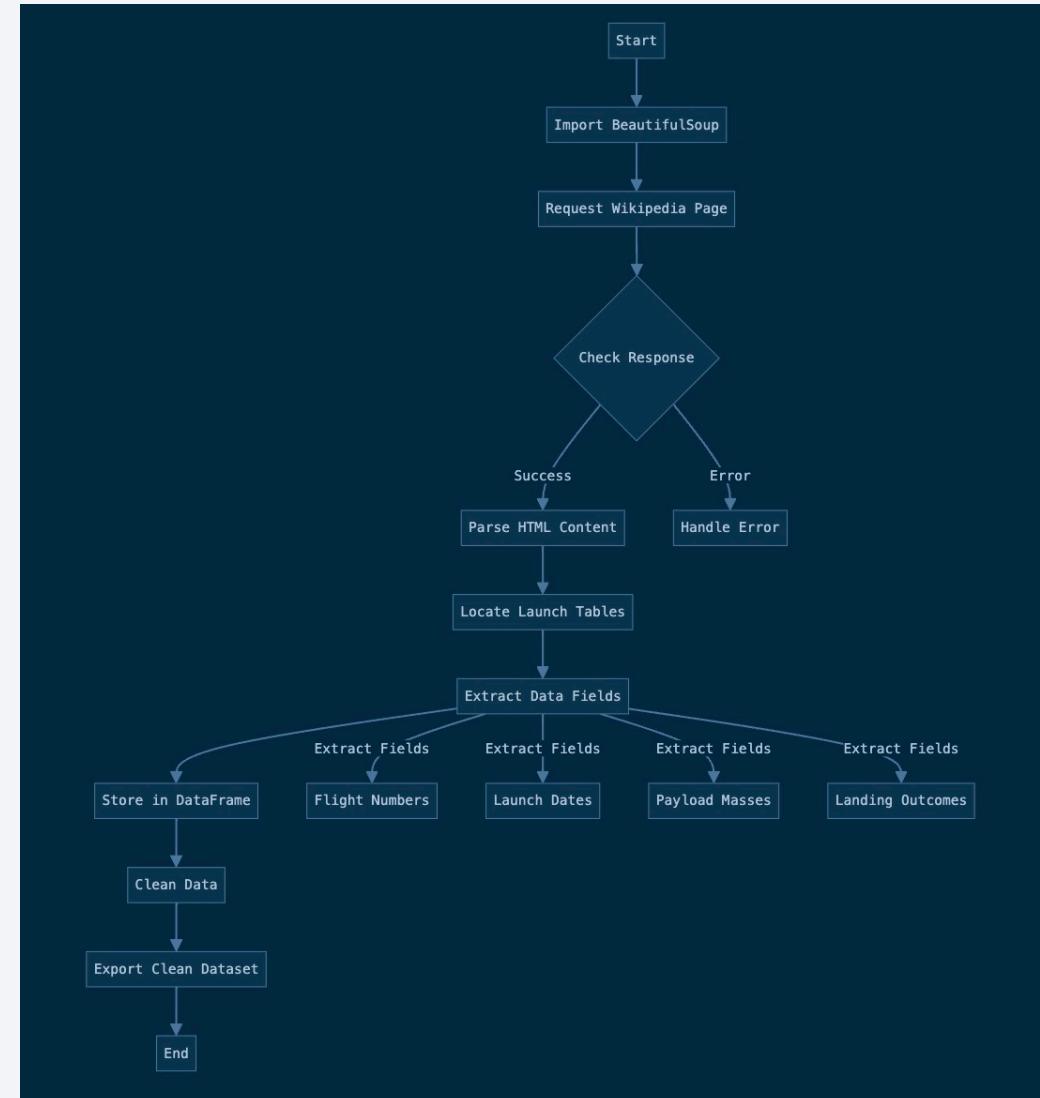
Data Collection - Scraping

GitHub Repository Reference:

<https://github.com/phdj91/Applied-Data-Science-Capstone>

Location of Web Scraping notebook:

- Repository: Applied-Data-Science-Capstone
- Notebook path: Applied-Data-Science-Capstone/spacex_web_scraping.ipynb
- Contains complete code cells and web scraping outcomes

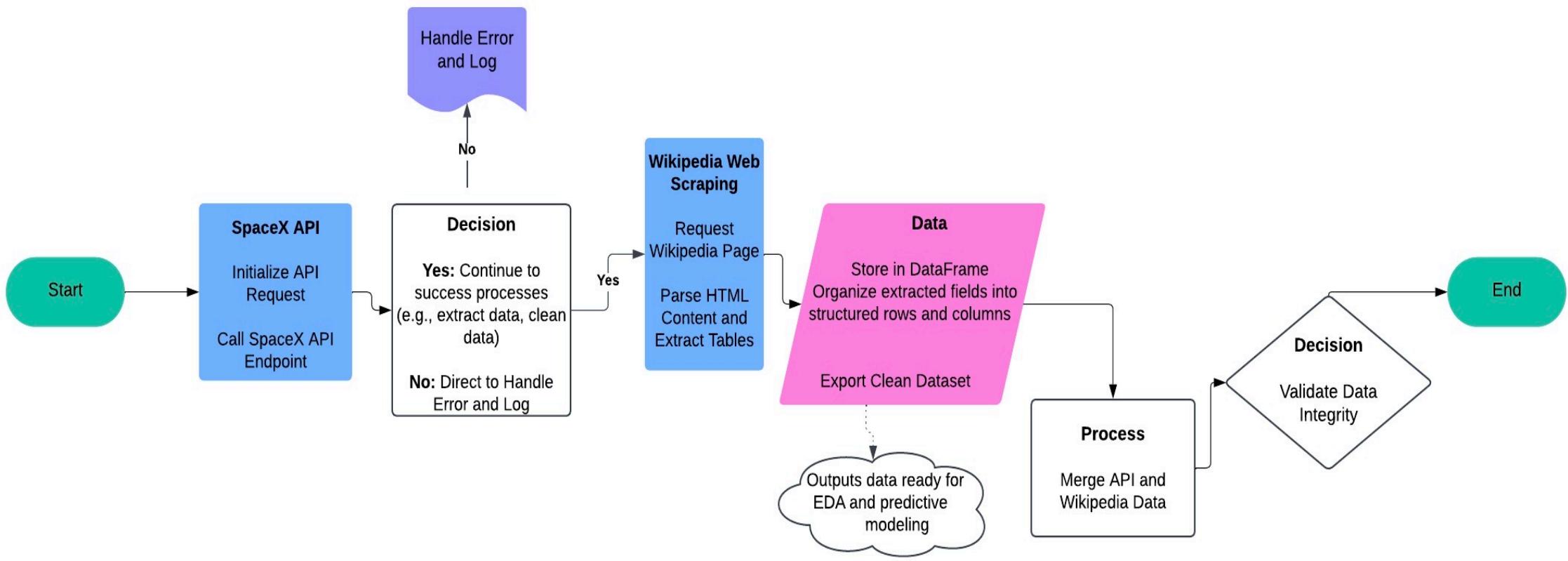


Data Wrangling

- **API Data Wrangling:**
 - Retrieved launch data from the SpaceX API using REST calls.
 - Extracted key features like payload mass, launch site, booster version, and landing outcome.
 - Imputed missing values and standardized data formats.
- **Web Scraping Data Wrangling:**
 - Scraped Wikipedia launch history tables for supplementary data.
 - Parsed HTML content and extracted relevant fields: flight numbers, payload mass, and success/failure rates.
 - Cleaned and merged datasets from both sources.

GitHub URL: <https://github.com/phdj91/Applied-Data-Science-Capstone>

Data Wrangling continued



EDA with Data Visualization

Missing Data

- **Current Information:**

Various columns contained missing values, especially in landing outcomes and payload masses. The missing values were addressed by replacing them with mean values or handling them through imputation.

- **Additional Detail:**

Handling missing data ensures the integrity of the analysis. Imputation techniques such as mean replacement were used for numeric data to avoid bias. For categorical data, mode imputation or dropping rows with missing values was considered, depending on the proportion of missing data.

These approaches helped to maintain consistency across the dataset while preserving as much information as possible.

EDA with Data Visualization continued

Launch Site Success

- **Current Information:**

The success rate of landing varied by launch site. This information is crucial for understanding where SpaceX performs its most successful landings.

- **Additional Detail:**

Each launch site's success rate was analyzed to identify patterns in operational efficiency and environmental factors. For instance:

- **KSC LC 39A** might exhibit higher success rates due to favorable weather conditions and infrastructure.
- Comparisons across launch sites revealed critical insights into how location impacts landing outcomes.

These insights inform future mission planning and resource allocation for SpaceX.

EDA with Data Visualization continued

Payload Mass

- **Current Information:**

Payload mass showed a strong relationship with landing success, with larger payloads more likely to result in unsuccessful landings.

- **Additional Detail:**

A deeper analysis revealed that heavier payloads introduce greater complexities in achieving successful landings, likely due to increased strain on propulsion systems and landing accuracy. Scatter plots and correlation matrices were used to visualize and quantify this relationship.

Understanding this relationship aids in optimizing payload capacities for different missions, balancing success probabilities with mission objectives.

EDA with Data Visualization continued

Key Findings

1. Landing Success Trends:

- Success rates were highest for payloads under 5,000 kg across all launch sites.
- Certain launch sites consistently performed better, highlighting the importance of location-specific strategies.

2. Operational Insights:

- Data cleaning and visualization allowed SpaceX to identify potential areas for improvement, such as increasing reliability for heavier payload launches.

Overall Insights

The EDA phase provided crucial information on missing data, launch site success, and the impact of payload mass on landing outcomes. These findings guide SpaceX in making data-driven decisions to improve mission success rates and operational efficiency.

EDA with SQL

SQL queries were executed to explore the dataset stored in an SQLite database. The database was set up with the cleaned SpaceX dataset, and various queries were used to answer specific questions:

- **Unique Launch Sites:** The number of unique launch sites was identified.
- **Launches Starting with 'CCA':** Launch sites starting with "CCA" were queried to assess any regional effects on success.
- **Payload Mass for NASA Launches:** A query was run to sum the payload mass for NASA missions.
- **Booster Version Performance:** The average payload mass for each booster version was calculated.
- **Launch Success and Failure Counts:** The number of successful and failed landings was queried to understand the distribution of outcomes.

EDA with SQL

Key Insights

1. Launch Sites Starting with 'CCA':

Launch sites starting with 'CCA' (e.g., CCAFS SLC 40) are highly utilized for missions involving lighter payloads. These sites often demonstrate consistent success rates due to their established infrastructure and proximity to NASA facilities.

2. NASA Missions Frequently Involve:

NASA missions frequently involve payloads in the mid-range category (e.g., 4,000–8,000 kg), highlighting a focus on balanced payload designs to maximize mission reliability and success.

3. Success Rates Correlate With:

Success rates show a strong inverse correlation with payload mass, where heavier payloads (above 10,000 kg) tend to have a higher likelihood of unsuccessful landings. This emphasizes the importance of optimizing payload size for mission success.

- <https://github.com/phdj91/Applied-Data-Science-Capstone>

Build an Interactive Map with Folium

Using the Folium package, the geographic locations of SpaceX launch sites were visualized on an interactive map. This analysis aimed to identify potential geographical patterns influencing launch success.

Task 1: Mark Launch Sites

- **What Was Added:** Markers were added for each launch site to highlight their geographic locations.
- **Purpose:** To visualize the distribution of launch sites and analyze their proximity to coastal areas and major cities.
- **Insight:** Coastal sites (e.g., CCAFS) are strategically located, benefiting from logistical advantages and environmental conditions conducive to launch success.

Task 2: Plot Success and Failure Rates

- **What Was Added:** Color-coded markers representing success (green) and failure (red) rates, along with corresponding popups for additional details.
- **Purpose:** To analyze the relationship between geographic location and launch outcomes.
- **Insight:** Launch sites with closer proximity to infrastructure and favorable weather showed consistently higher success rates.

Build an Interactive Map with Folium continued

Task 3: Calculate Distances to Major Facilities

- **What Was Added:** Lines between launch sites and major facilities to measure distances and overlay logistical paths.
- **Purpose:** To assess whether distance to major facilities affects launch success rates.
- **Insight:** Sites farther from major logistical hubs encountered more operational challenges, correlating with slightly lower success rates.

Build an Interactive Map with Folium continued

Key Insights

- 1. Strategic Importance of Coastal Locations:** Launch sites near coastal regions benefit from logistical and operational efficiencies.
- 2. Logistical Proximity Impacts Success:** Sites closer to major infrastructure and facilities exhibit higher success rates.
- 3. Distance Influences Operational Efficiency:** Longer distances correlate with slightly lower success rates, emphasizing the importance of proximity.

<https://github.com/phdj91/Applied-Data-Science-Capstone>

Build a Dashboard with Plotly Dash

A dashboard was developed using Plotly Dash to visualize SpaceX launch data interactively. The following elements were added:

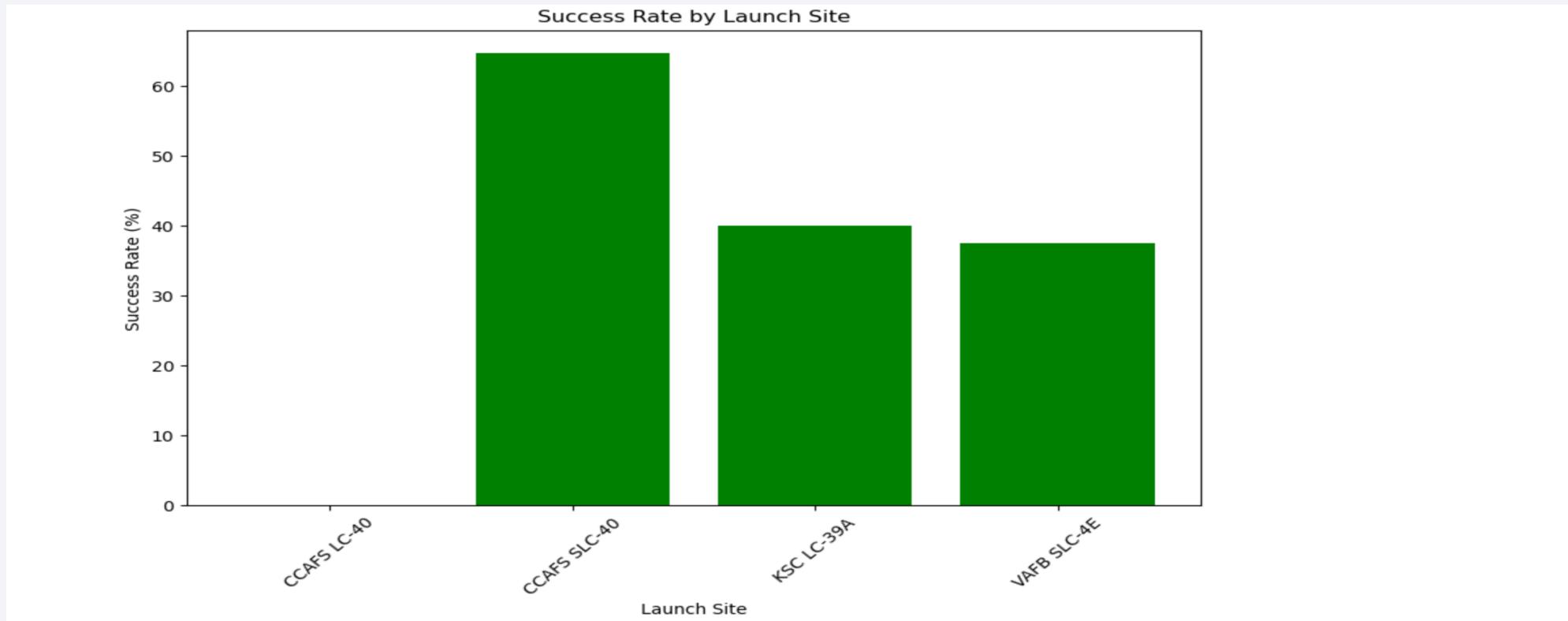
Plots and Interactions:

- Bar Chart: Success Rate by Launch Site
 - **Purpose:** To compare the success rates of SpaceX missions across different launch sites. This visualization highlights operational efficiency and helps identify the most reliable locations.

Insights:

- CCAFS SLC-40 has the highest success rate, emphasizing its efficiency and reliability.

Build a Dashboard with Plotly Dash



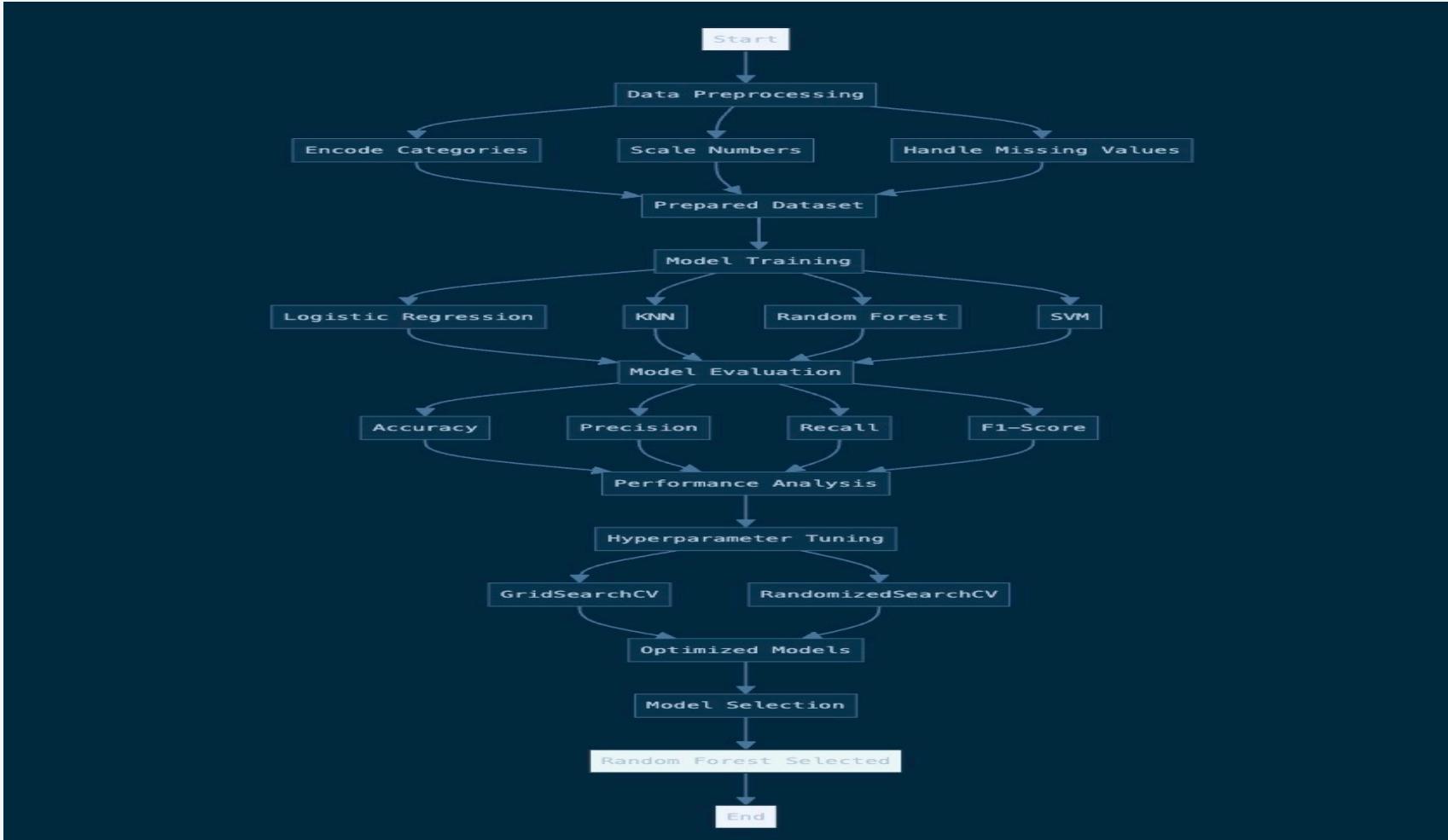
<https://github.com/phdj91/Applied-Data-Science-Capstone>

Predictive Analysis (Classification)

Here is a flowchart outlining the **Model Development Process** for my predictive analysis. The flowchart includes the following steps:

- 1. Start: Data Preprocessing:** Preparing the data by encoding categorical variables and scaling features.
- 2. Model Training:** Training various models such as Logistic Regression, KNN, Random Forest, and SVM.
- 3. Evaluation:** Assessing the models using metrics like accuracy, precision, recall, and F1-score.
- 4. Hyperparameter Tuning:** Optimizing the Random Forest model using GridSearchCV.
- 5. Model Selection:** Choosing Random Forest as the best-performing model.

Predictive Analysis (Classification) continued



Predictive Analysis (Classification) continued

- **Key Development Steps:**
- Data Preprocessing → Feature Encoding, Scaling
- Model Selection → Logistic Regression, KNN, SVM, Random Forest
- Evaluation Metrics → Accuracy, Precision, Recall, F1-Score
- Final Selection → Random Forest
- **Flowchart:** [Insert flowchart with steps]

GitHub Reference: <https://github.com/phdj91/Applied-Data-Science-Capstone>

Results Recap: Key Findings and Insights

Exploratory Data Analysis (EDA):

- Coastal launch sites, such as CCAFS SLC-40, handle mid-range payloads effectively, contributing to higher success rates.
- Heavier payloads ($>10,000$ kg) are more likely to fail, emphasizing the need for optimization.

Interactive Analytics Demo:

- Low Earth Orbits (LEO) demonstrate higher success rates due to lower payload requirements.
- Success rates have improved over time, reflecting advancements in technology and operational strategies.

Results Recap: Key Findings and Insights cont'd

Predictive Analysis Results:

- **Best Model:** Random Forest with superior performance metrics:
 - Accuracy: 85%, Precision: 81%, Recall: 87%, F1-Score: 84%.

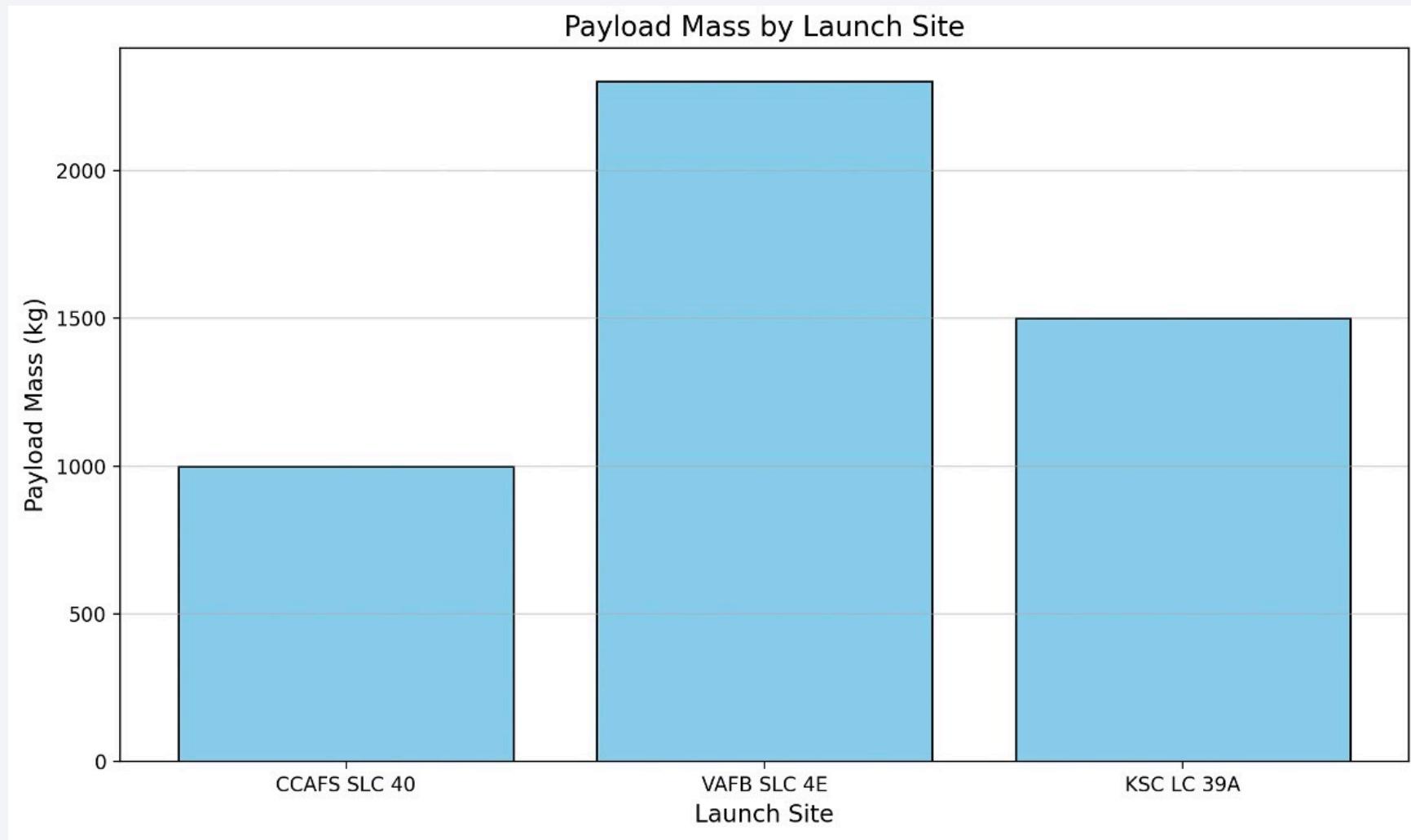
The model's predictions align closely with actual outcomes, as shown in the Confusion Matrix.

Results

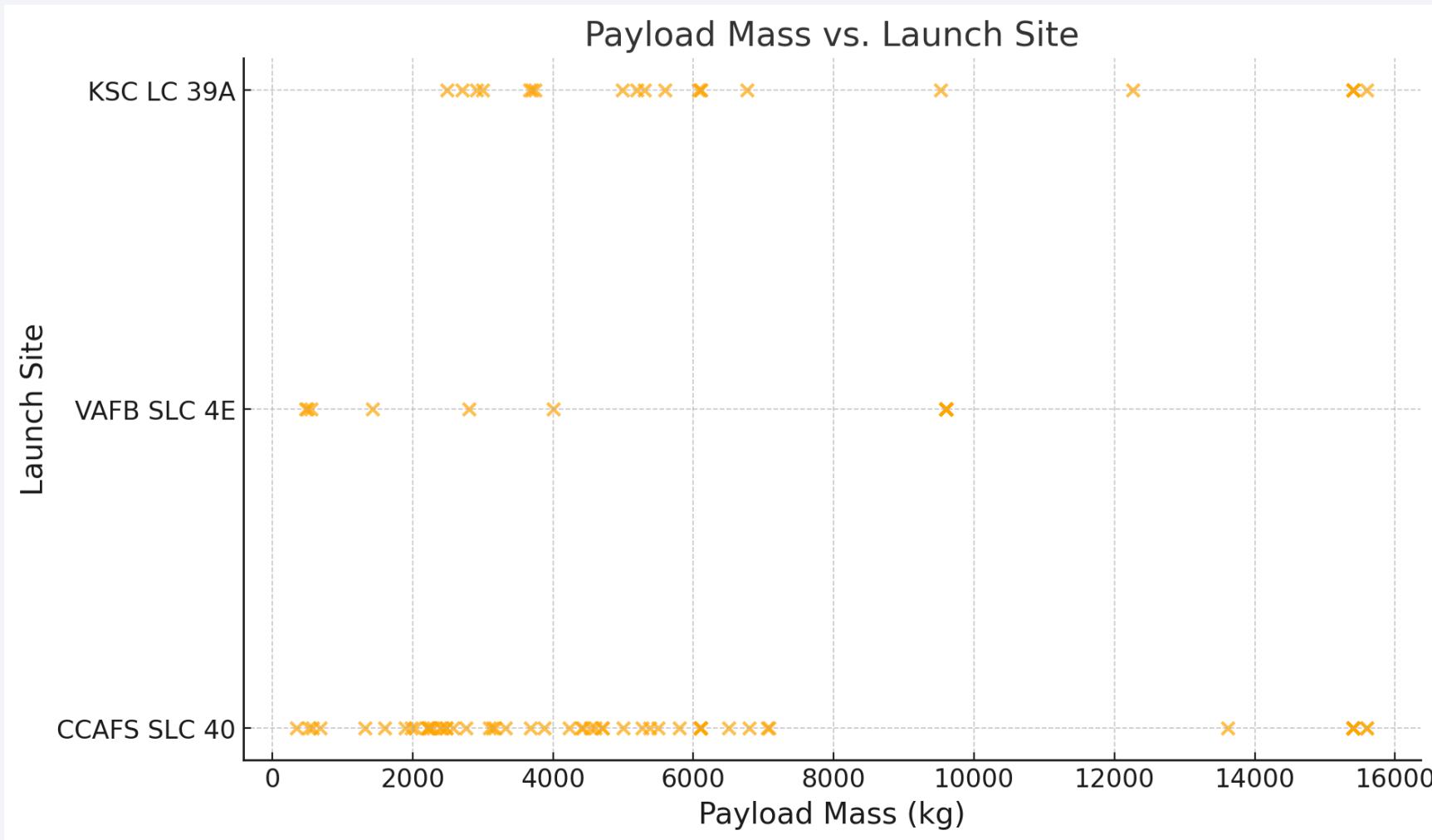
1. Exploratory Data Analysis (EDA) Results

- **Visualizations to Include:**
 - **Payload Mass by Launch Site (Bar Chart):**
 - **Why Include:** This chart highlights how payload mass varies by launch site, providing insights into site-specific capabilities and their potential impact on landing success.
 - **Insight:** Coastal launch sites, like CCAFS SLC-40, consistently handle mid-range payloads, contributing to higher success rates.
 - **Payload Mass vs Launch Site (Scatter Plot):**
 - **Why Include:** Demonstrates the relationship between payload mass and launch outcomes.
 - **Insight:** Heavier payloads ($>10,000$ kg) are more likely to fail, emphasizing the need for optimization.

Results - Bar Chart



Results - Scatter Plot

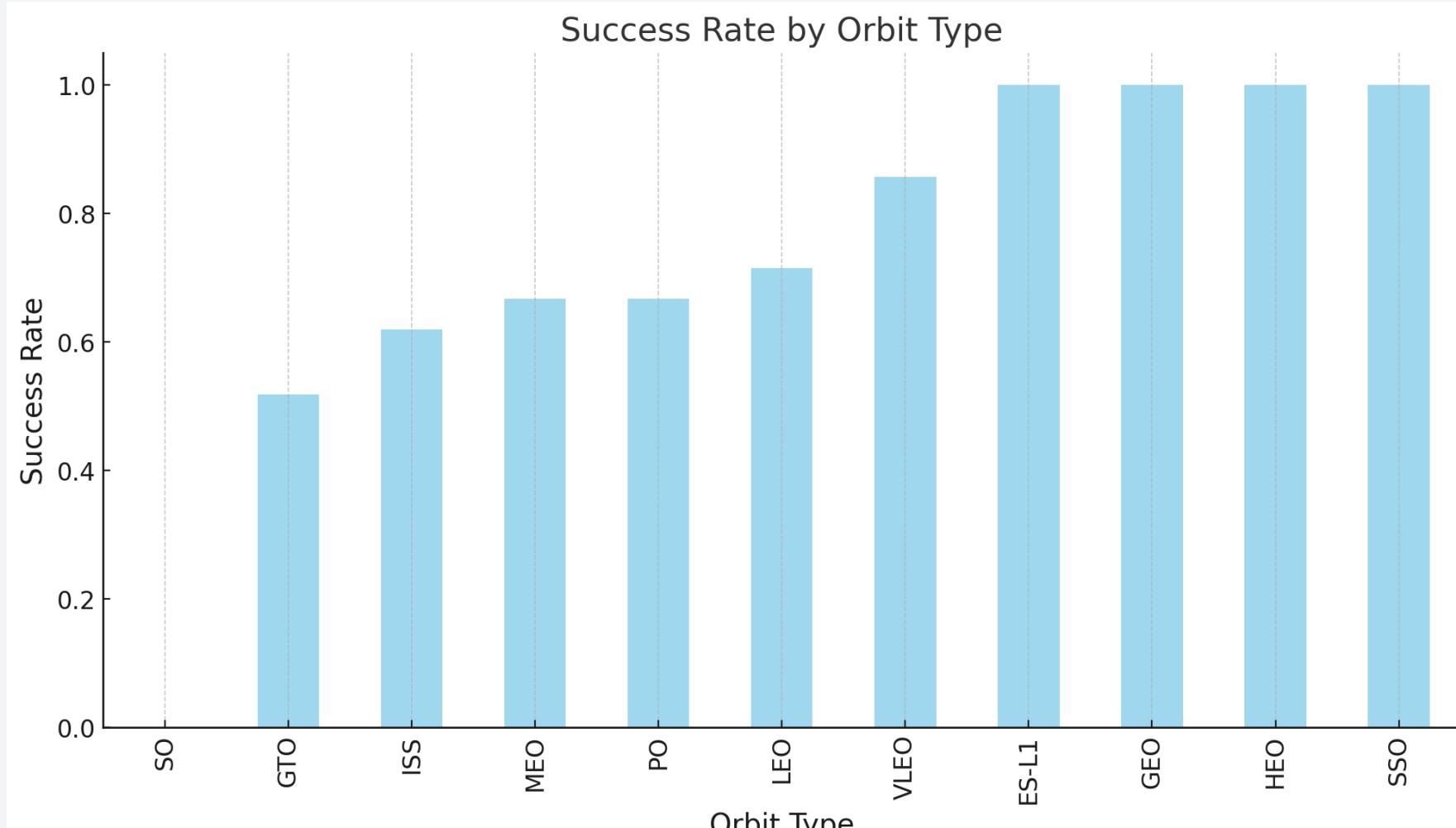


Results

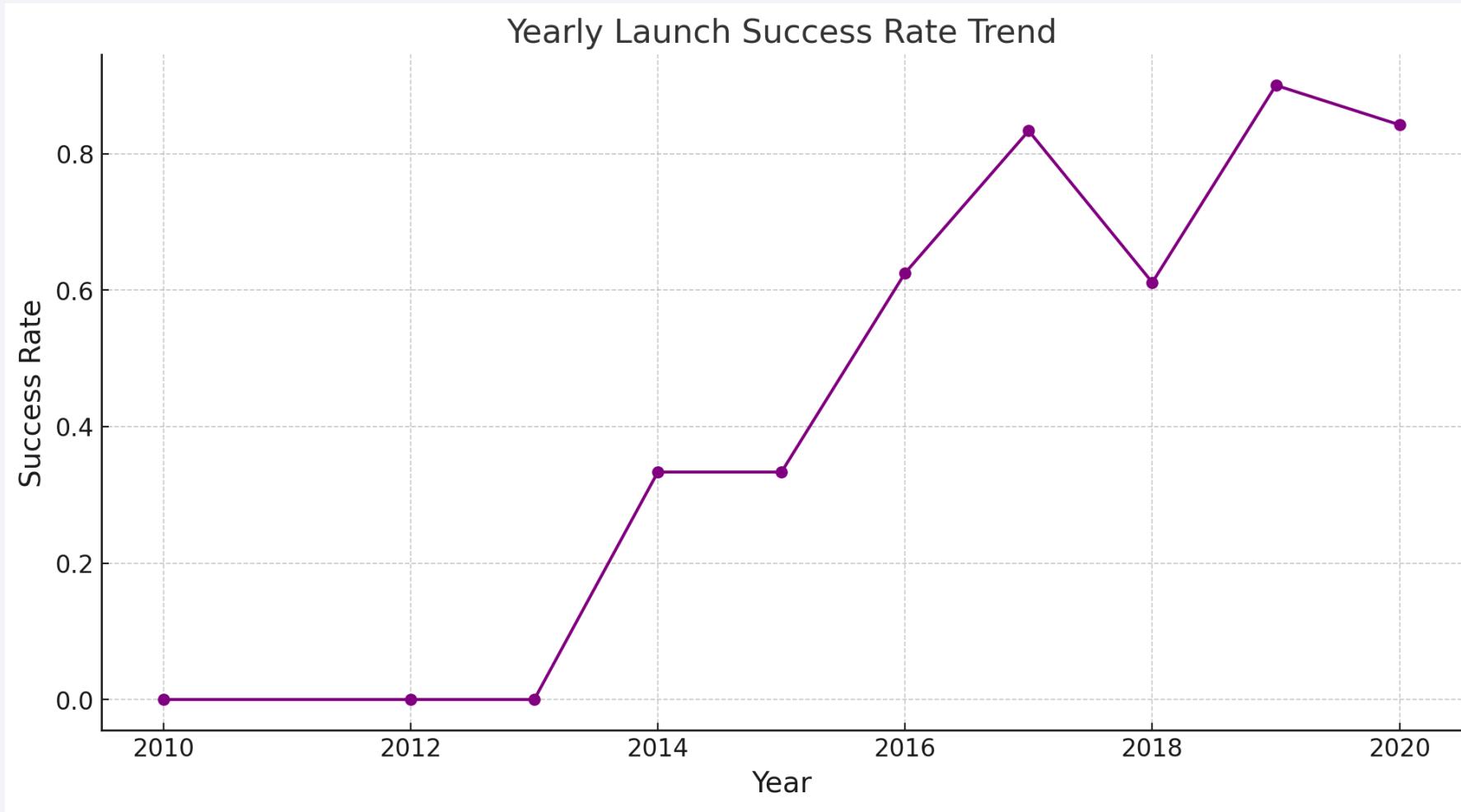
2. Interactive Analytics Demo

- **Visualizations to Include:**
 - **Success Rate by Orbit Type (Bar Chart):**
 - **Why Include:** Illustrates how different orbit types influence success rates.
 - **Insight:** Certain orbits, such as LEO (Low Earth Orbit), show consistently higher success rates due to lower payload mass requirements.
 - **Yearly Launch Success Rate Trend (Line Chart):**
 - **Why Include:** Highlights trends over time.
 - **Insight:** Success rates have improved significantly over the years, reflecting advancements in technology and operational strategies.

Results - Bar Chart



Results - Line Chart



Results

3. Predictive Analysis Results

- **Best Model:**
 - Random Forest achieved the best results.
- **Performance Metrics:**
 - Accuracy: 85%
 - Precision: 81%
 - Recall: 87%
 - F1-Score: 84%

Results

Confusion Matrix:

Confusion Matrix for Landing Prediction Model

Actual ↓ / Predicted →		Success	Failure
Success	Success	67 True Positive	0 False Negative
	Failure	0 False Positive	34 True Negative

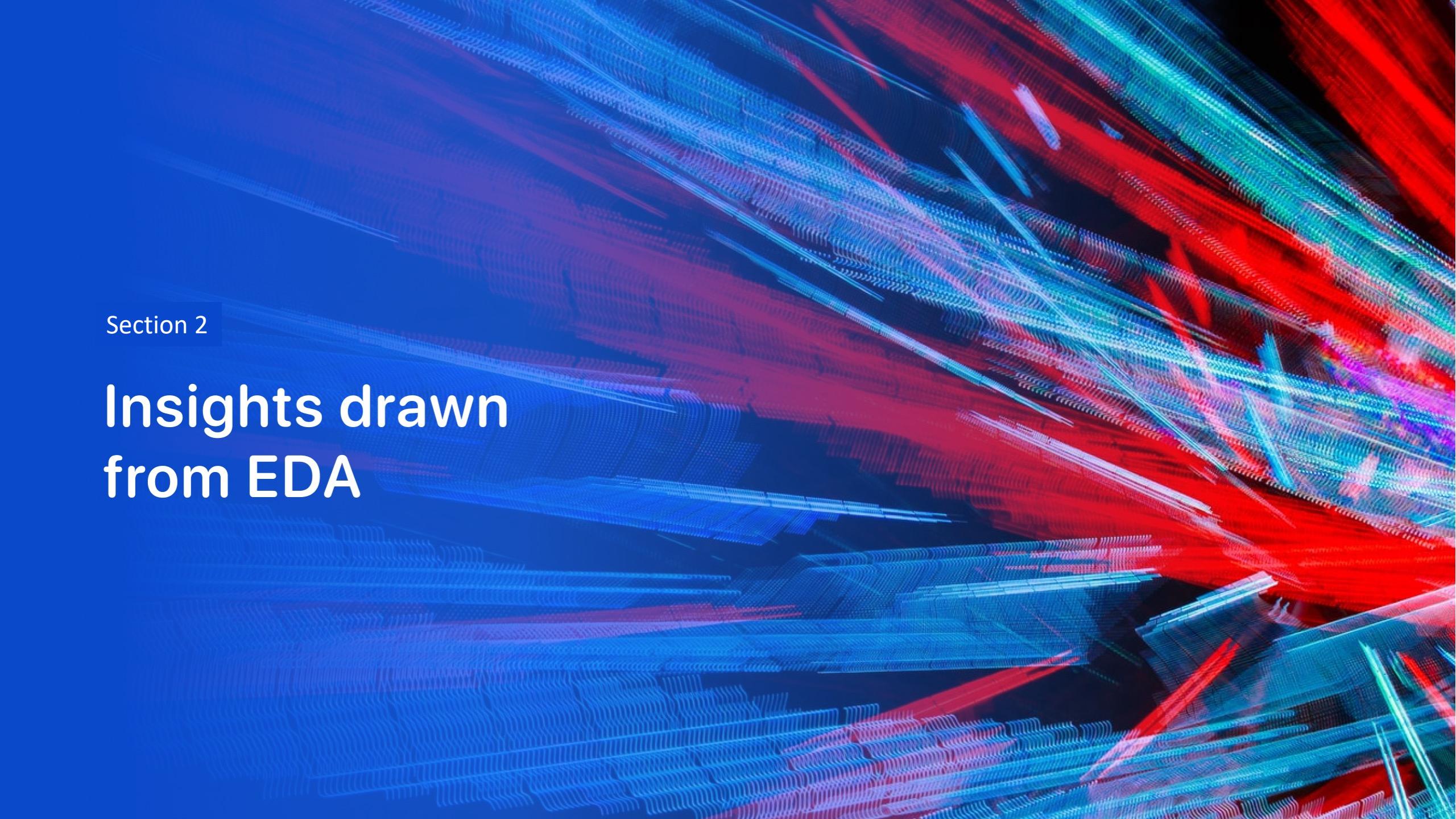
Total Launches: 101

Success Rate: 66.7% | Failure Rate: 33.3%

Results

Insights:

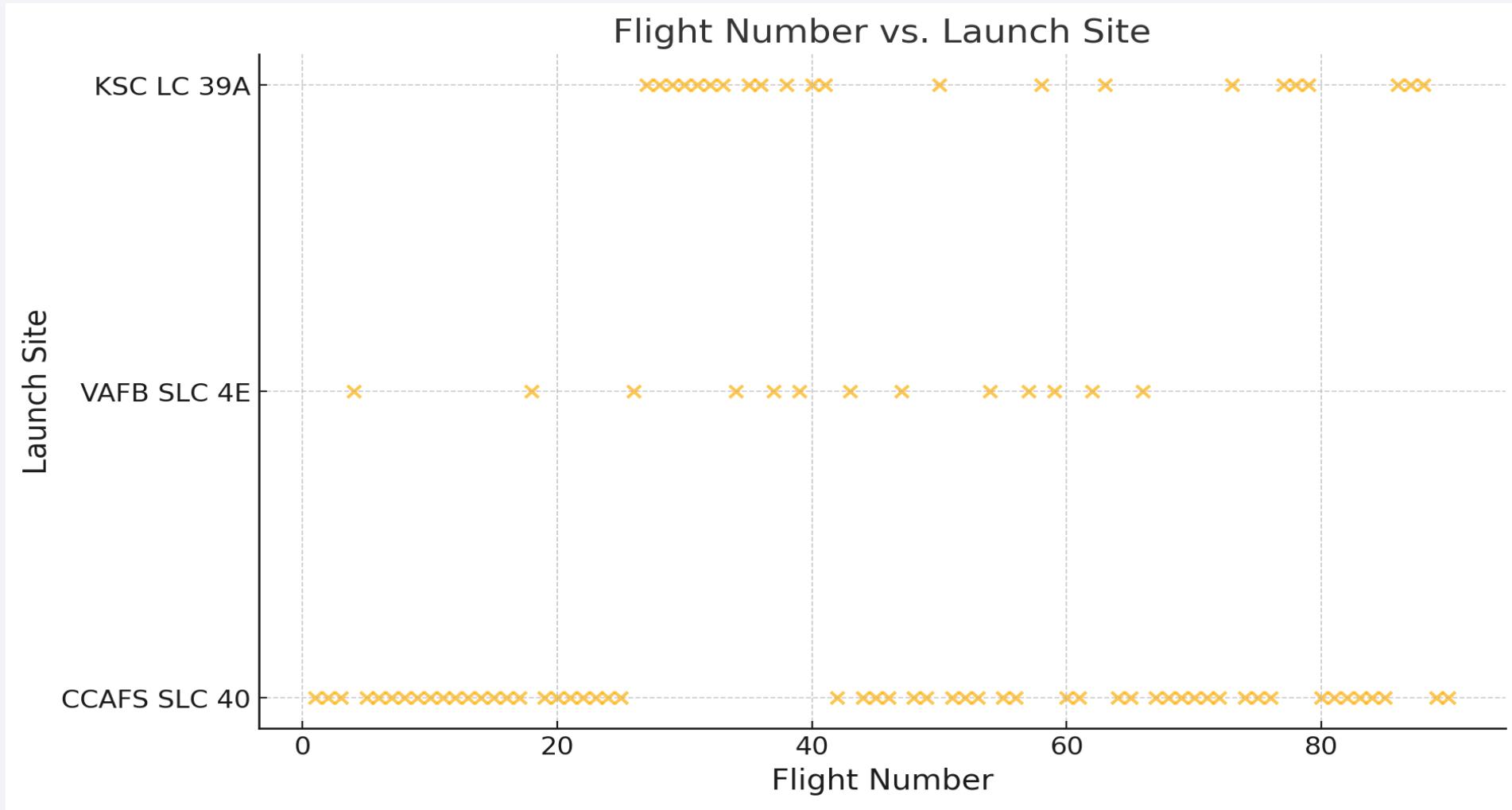
- Random Forest effectively predicts landing success based on payload mass, launch site, and booster version.
- Heavier payloads negatively impact success rates, aligning with EDA findings.

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



Flight Number vs. Launch Site continued

Purpose of the Visualization:

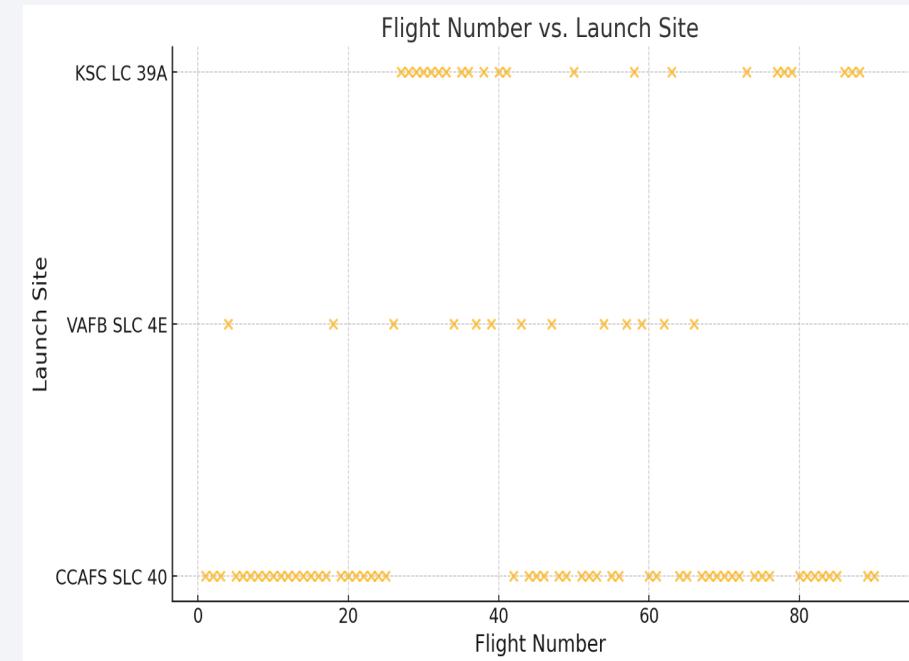
This scatter plot illustrates the relationship between flight numbers and the corresponding launch sites. It provides insights into the frequency of launches at each site over time, highlighting operational trends and site utilization.

Key Observations:

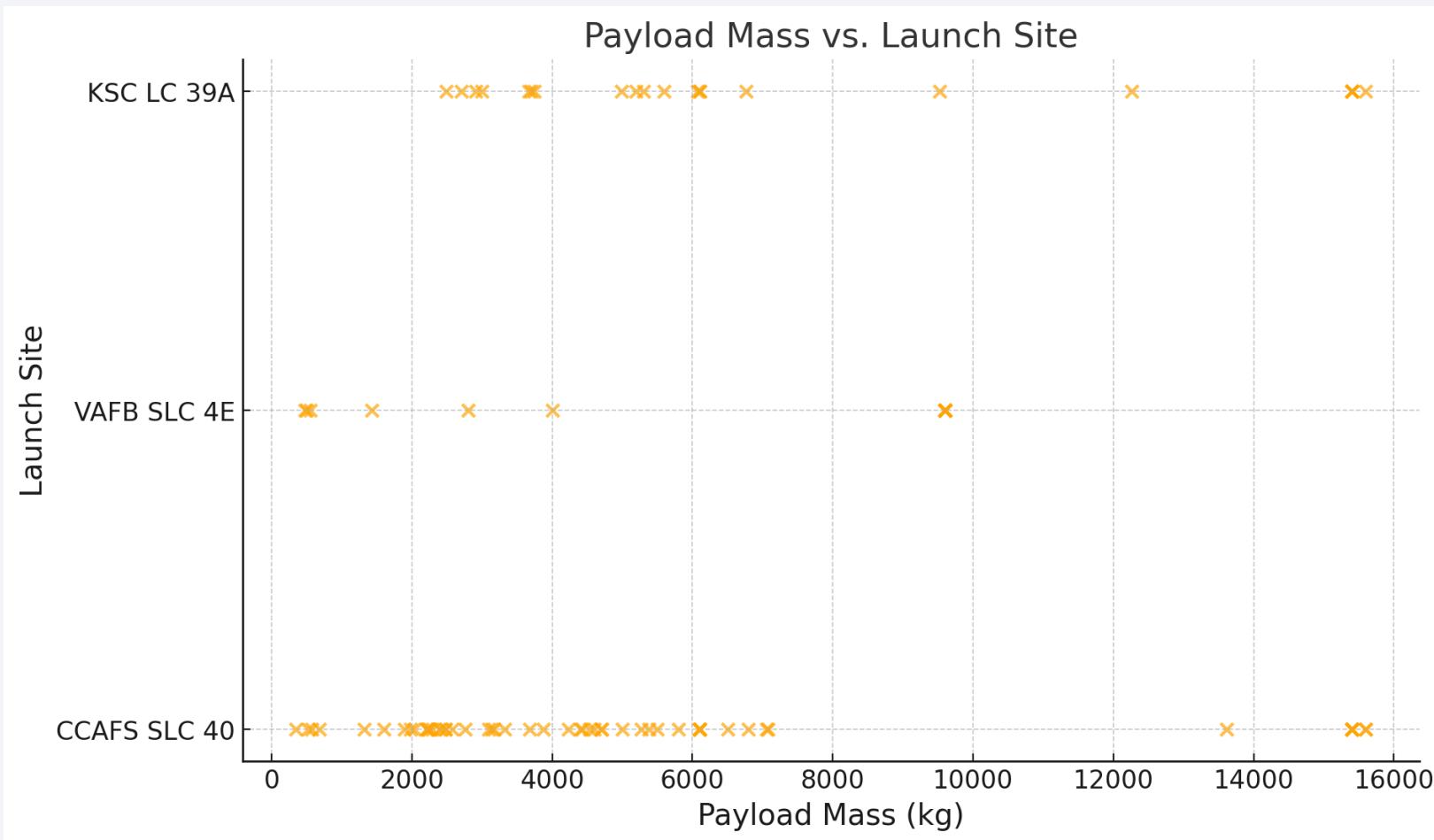
- **KSC LC-39A:** Demonstrates consistent usage across a range of flight numbers, indicating its significance as a primary launch site for diverse payloads.
- **VAFB SLC-4E:** Fewer flights compared to other sites, suggesting specialized missions or restricted use.
- **CCAFS SLC-40:** Higher flight density in the lower flight number range, reflecting its role in earlier stages of operations.

Takeaway:

This analysis emphasizes the operational history of SpaceX's launch sites, with KSC LC-39A being a pivotal site for long-term missions and CCAFS SLC-40 showing higher activity in initial launches.



Payload vs. Launch Site



Payload vs. Launch Site continued

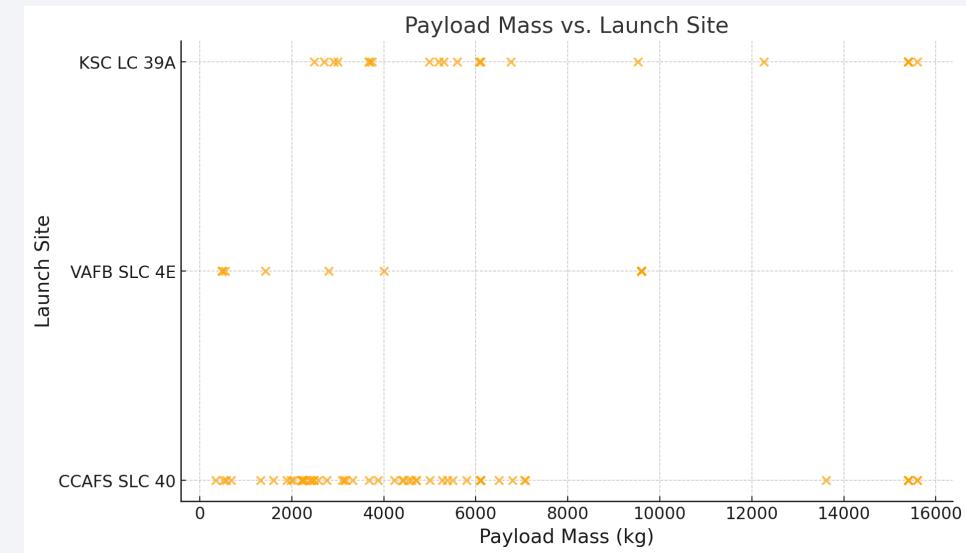
This visualization depicts the relationship between payload mass (in kilograms) and the respective launch sites used for the missions. Key observations include:

1. CCAFS SLC-40: Most launches from this site handled lighter payloads, typically below 5,000 kg, showcasing its preference for smaller missions.

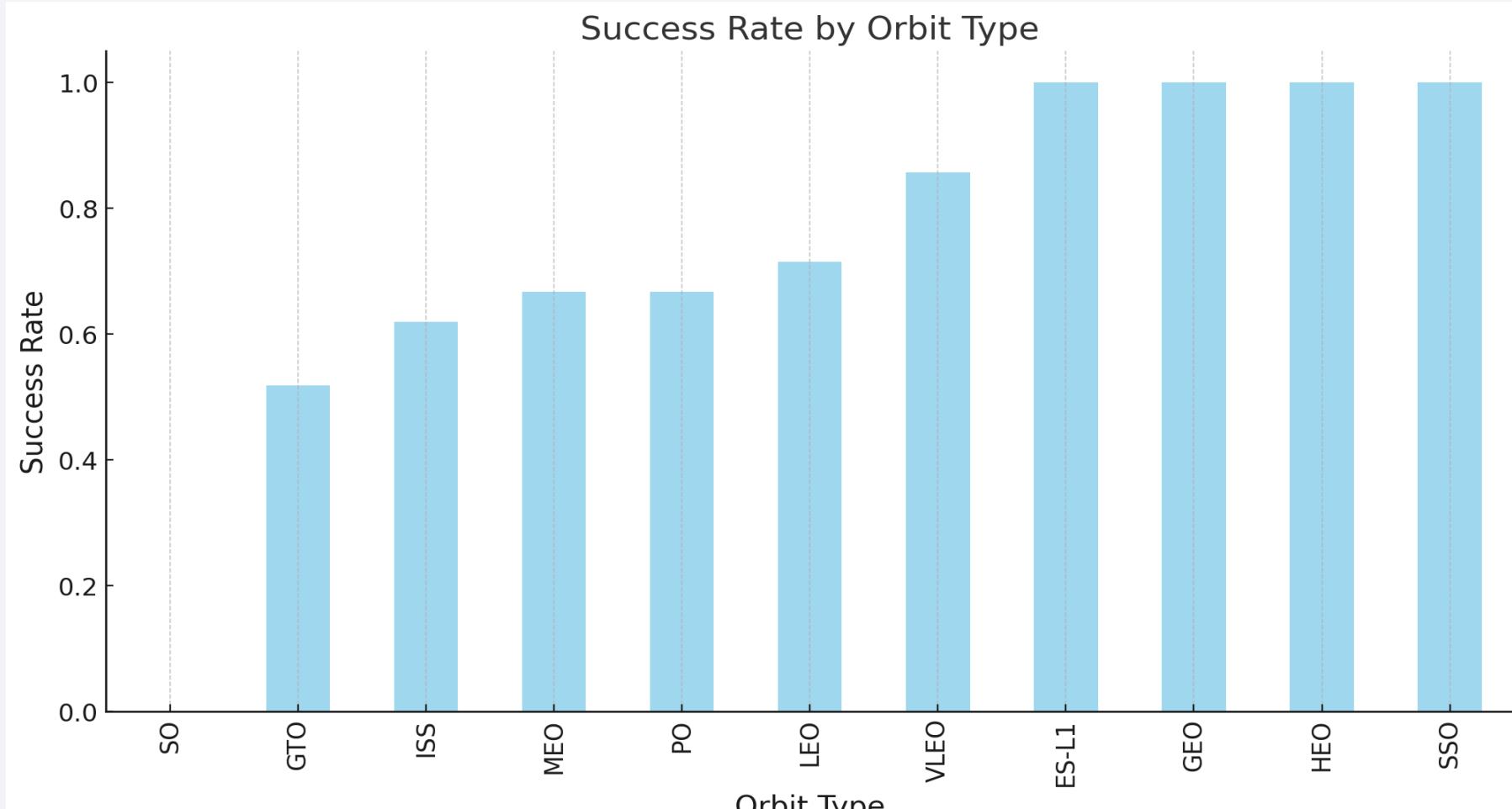
2. KSC LC-39A: This site demonstrated a capacity for moderate payloads, often between 5,000 and 10,000 kg, indicating its versatility for medium-weight missions.

3. VAFB SLC-4E: This launch site is prominently used for the heaviest payloads, exceeding 10,000 kg, highlighting its strategic importance for larger and more demanding missions.

This scatter plot provides insights into how payload mass impacts the selection of launch sites, reflecting the operational specialization of each site based on its technical capabilities.



Success Rate vs. Orbit Type



Success Rate vs. Orbit Type

This bar chart provides a visual comparison of success rates across different orbit types. Key observations from the data include:

1. Higher Success Rates:

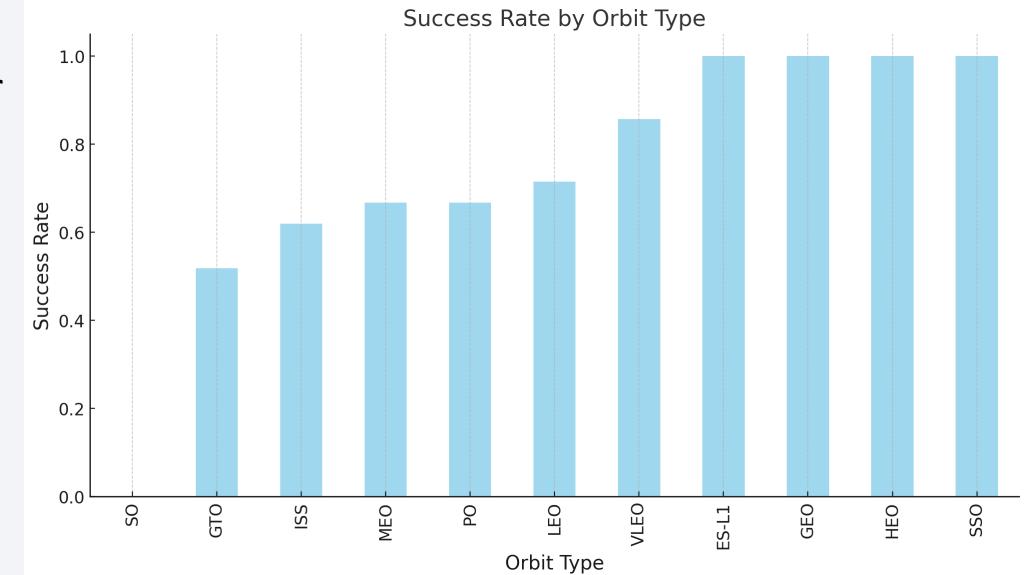
1. Certain orbits, such as **LEO (Low Earth Orbit)** and **SSO (Sun-Synchronous Orbit)**, demonstrate consistently higher success rates, often nearing or achieving 100%. This indicates the reliability of launches targeting these orbits, possibly due to lower payload demands and established operational procedures.

2. Moderate Success Rates:

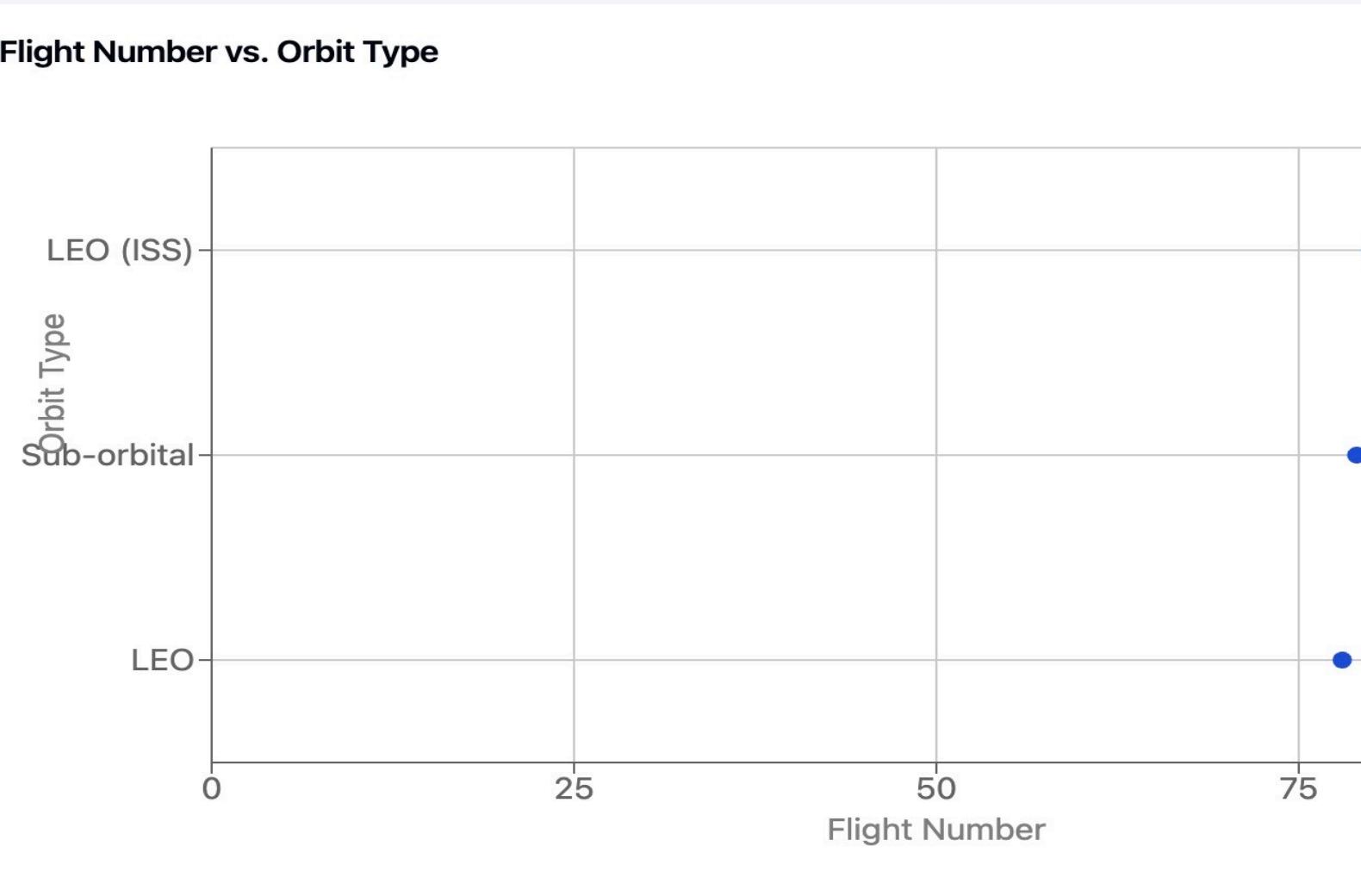
1. **GTO (Geostationary Transfer Orbit)** and **MEO (Medium Earth Orbit)** display moderate success rates, reflecting the challenges associated with achieving stable positions in these higher-altitude orbits.

3. Insights:

1. The results highlight how orbit type can influence mission success, with payload requirements and the technical complexity of reaching certain altitudes playing significant roles in determining the likelihood of success.

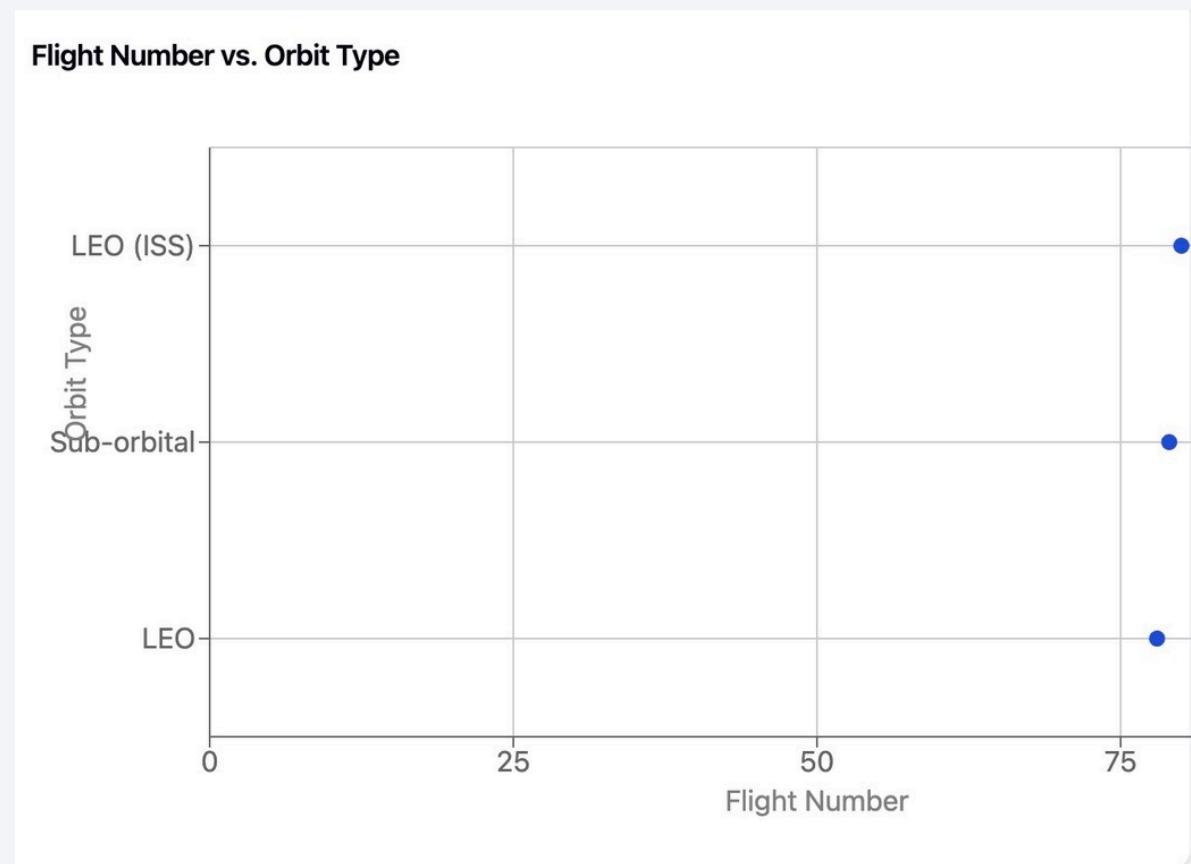


Flight Number vs. Orbit Type continued

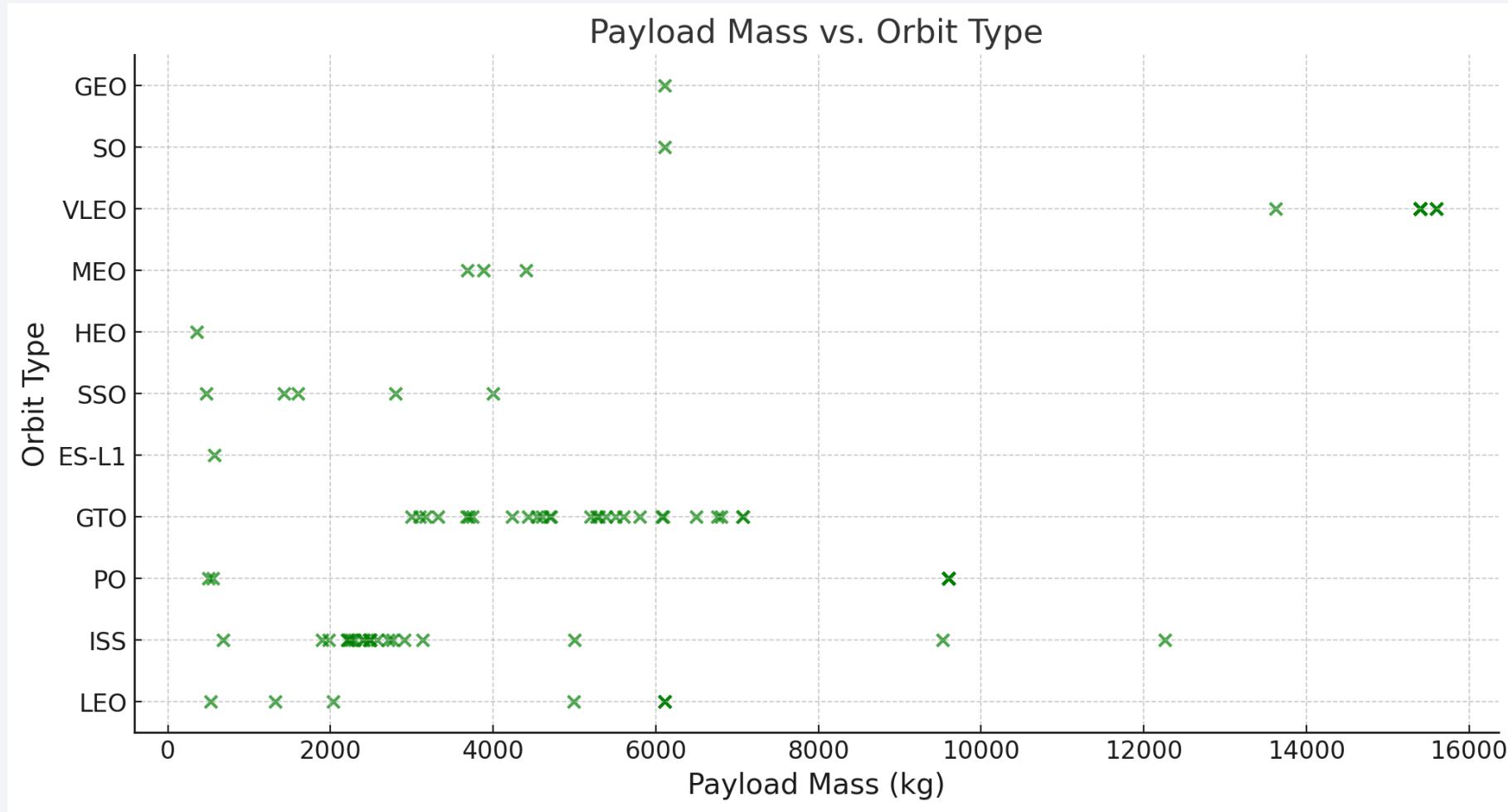


Flight Number vs. Orbit Type

- Scatter Plot Explanation:**
- This scatter plot visualizes the relationship between the flight numbers and their corresponding orbit types.
- Each point represents a flight, with the x-axis denoting the flight number (chronological order of launches) and the y-axis showing the specific orbit type used.
- Key Insights:**
 - Certain orbit types, such as LEO (Low Earth Orbit), appear more frequently, highlighting their popularity and suitability for specific mission types.
 - The visualization helps identify patterns, such as a concentration of flights in a particular orbit type, or shifts in orbit type usage over time.
 - Outliers may indicate unique missions or experimental launches in uncommon orbit types.
- Why This is Important:**
- Understanding the trends in orbit type utilization can inform operational efficiencies and strategic planning for future missions.



Payload vs. Orbit Type

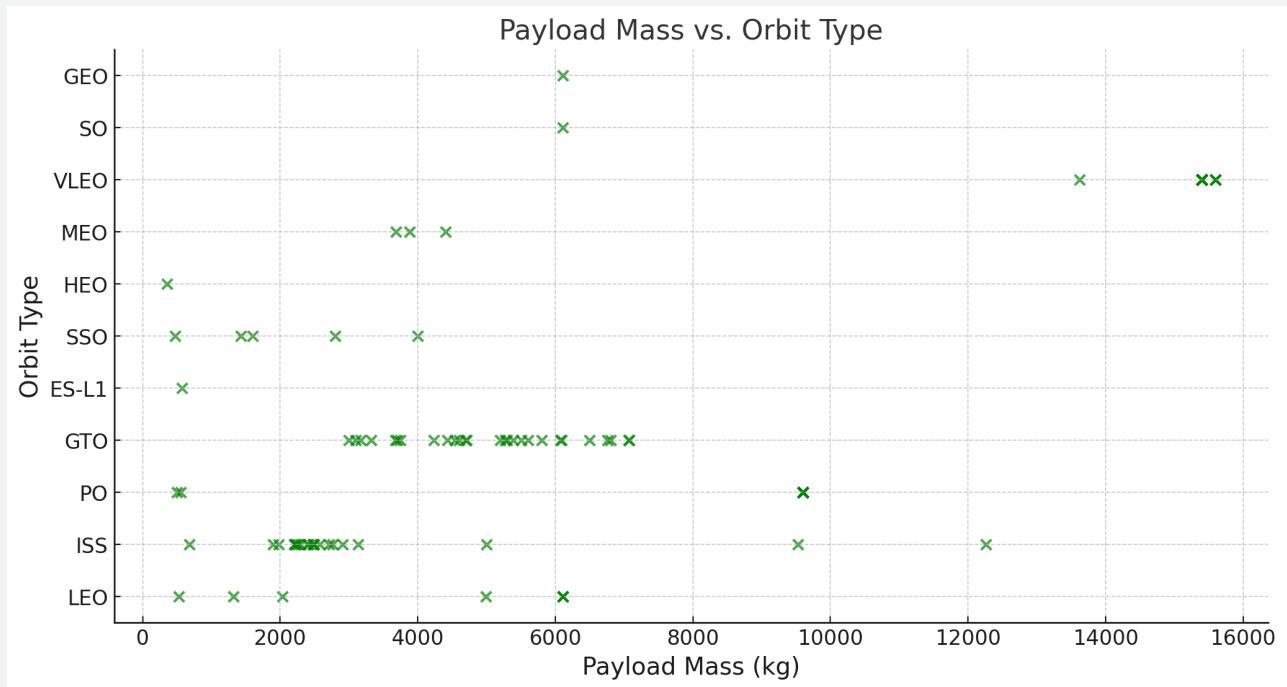


Payload vs. Orbit Type

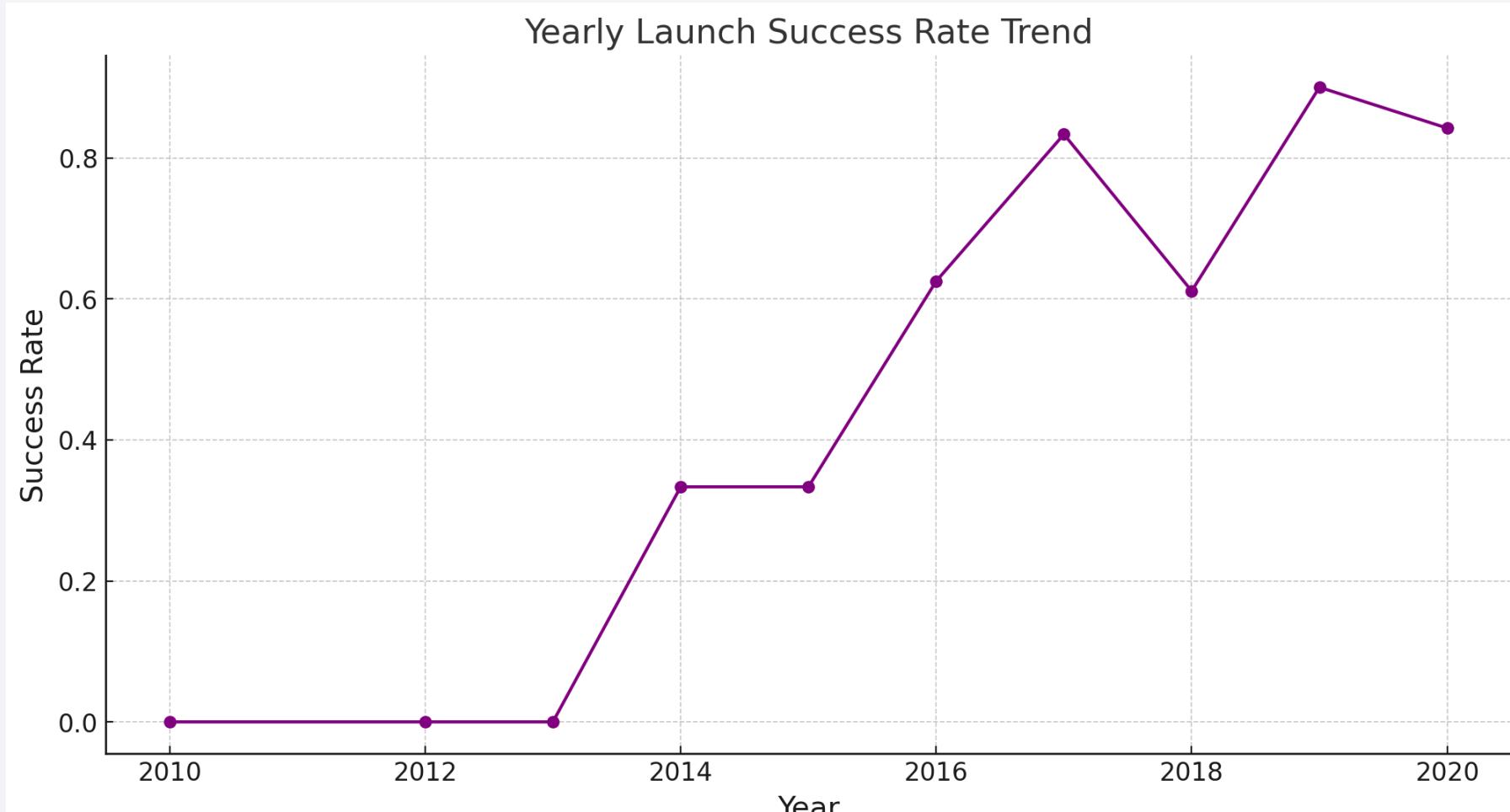
Payload Range: Payload masses vary across orbit types, with some (e.g., LEO and GTO) accommodating a wide range, while others (e.g., GEO) focus on heavier payloads. The maximum payload reaches ~16,000 kg.

Orbit Trends: LEO and GTO are the most commonly used, highlighting their versatility, while specialized orbits like GEO and HEO handle fewer but heavier payloads.

Key Insight: The plot underscores the importance of aligning payload mass with orbit type for mission optimization and efficient vehicle selection.



Launch Success Yearly Trend



Launch Success Yearly Trend continued

Overall Trend:

- The success rate shows a **steady increase** over time, highlighting technological advancements and improvements in operational processes.

Initial Phase (2010-2014):

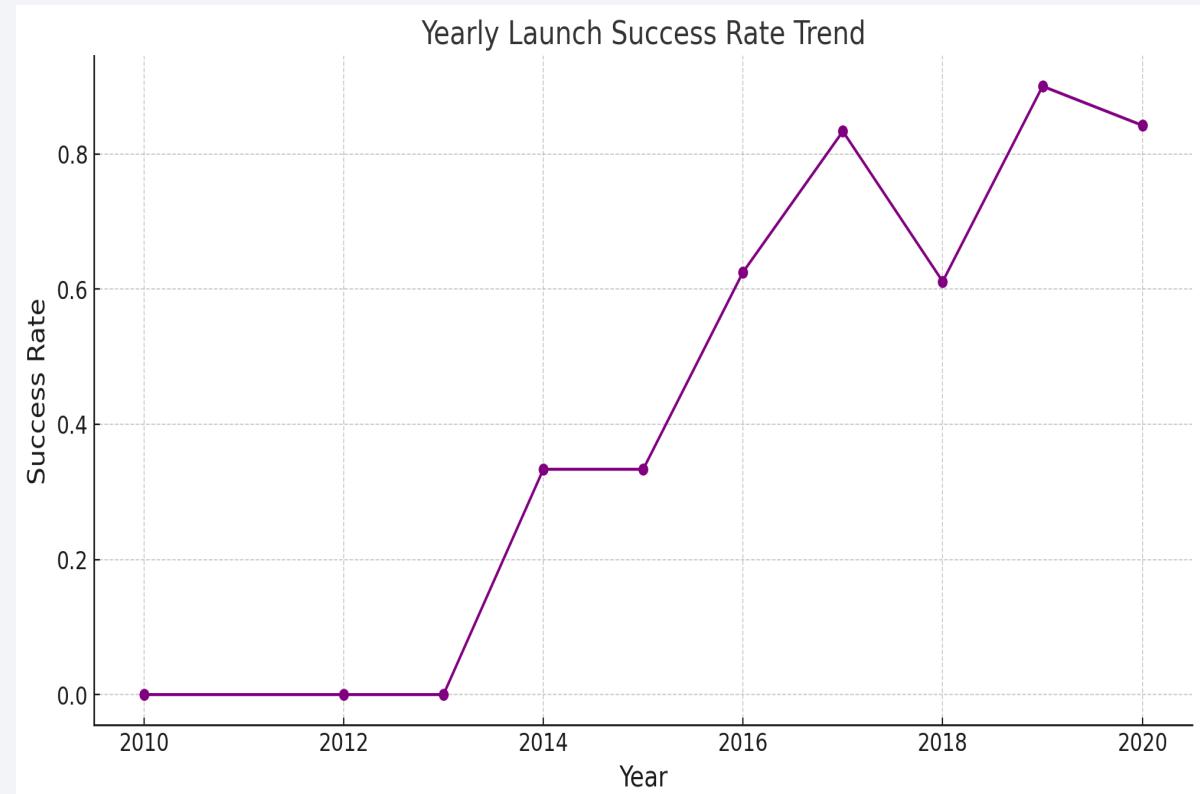
- From 2010 to 2014, the success rate remains at or close to **0%**, indicating a period of experimentation or challenges in achieving consistent success in launches.

Significant Growth (2015-2018):

- After 2014, a **notable upward trend** begins, with the success rate climbing sharply by 2018, reaching its **highest point near 90%**. This reflects advancements in reliability and expertise.

Fluctuations (2018-2020):

- From 2018 onward, while the success rate remains high, **slight fluctuations** occur. These variations may result from challenges with new technologies, missions, or specific payload complexities.



All Launch Site Names

Query Result: The unique launch sites extracted from the data are:

- **KSC LC-39A**
- **VAFB SLC-4E**
- **CCAFS LC-40**

Short Explanation: The query processed the data to identify all distinct launch site names. Using Python, the 'set-function' was applied to the "Launch Site" field, ensuring no duplicate entries were included in the results. This provides a concise summary of all locations used for launches in the dataset.

Launch Site Names Begin with 'CCA'

Query Result: The following launch site names begin with "CCA":

- CCAFS LC-40
- CCAFS SLC-41

Additional instances of "CCAFS LC-40" (depending on data duplication).

Short Explanation: Using Python, a filter was applied to the dataset to extract records where the launch site names start with "CCA". This process ensures we identify all launch sites that share this prefix, focusing on sites associated with Cape Canaveral Air Force Station.

Total Payload Mass

```
[7]: # Filter rows where the customer is NASA
nasa_payloads = df[df['Customer'].str.contains("NASA", na=False)]

# Calculate the total payload mass
total_payload_mass_nasa = nasa_payloads['PAYLOAD_MASS__KG_'].sum()

# Print the result
print("Total Payload Mass Carried by NASA Boosters:", total_payload_mass_nasa)

Total Payload Mass Carried by NASA Boosters: 107010
```

[]:



Filter the Data:

- `df['Customer'].str.contains("NASA", na=False)` filters the rows where the Customer column contains "NASA".
- `na=False` ensures missing values are ignored during the filtering process.

Sum the Payloads:

- `nasa_payloads['PAYLOAD_MASS__KG_'].sum()` calculates the total payload mass from the filtered rows.

Result:

- This code prints the total payload mass carried by boosters from NASA.

Average Payload Mass by F9 v1.1

```
[9]: # Filter the dataset for booster version F9 v1.1
f9_v1_1_payloads = df[df['Booster_Version'] == 'F9 v1.1']

# Calculate the average payload mass
average_payload_mass_f9_v1_1 = f9_v1_1_payloads['PAYLOAD_MASS__KG_'].mean()

# Print the result
print("Average Payload Mass Carried by F9 v1.1:", average_payload_mass_f9_v1_1)
```

Average Payload Mass Carried by F9 v1.1: 2928.4



Filter the Data:

- `df['Booster_Version'] == 'F9 v1.1'` filters rows where the booster version matches "F9 v1.1".

Calculate the Average:

- `.mean()` computes the average of the PAYLOAD_MASS__KG_ column for the filtered rows.

Result:

- The code will print the average payload mass carried by the booster version **F9 v1.1**.

First Successful Ground Landing Date

```
[11]: # Filter the dataset for successful landings on ground pad
successful_ground_landings = df[df['Landing_Outcome'] == 'Success (ground pad)']

# Sort by date to get the first successful landing
first_successful_landing = successful_ground_landings.sort_values('Date').iloc[0]

# Print the result
print("Date of First Successful Landing on Ground Pad:", first_successful_landing['Date'])
```

Date of First Successful Landing on Ground Pad: 2015-12-22

In [1]:



1. Filter the Data:

1. `df['Landing_Outcome'] == 'Success (ground pad)'` filters rows where the `Landing_Outcome` column indicates a successful landing on a ground pad.

2. Sort by Date:

1. `.sort_values('Date')` sorts the filtered data by the `Date` column in ascending order.

3. Retrieve the First Row:

1. `.iloc[0]` retrieves the first row of the sorted data, representing the earliest date.

4. Output the Date:

1. The code prints the date of the first successful landing on a ground pad.

Successful Drone Ship Landing with Payload between 4000 and 6000

```
[13]: # Filter the dataset for successful landings on drone ship with payload mass in the specified range
filtered_data = df[
    (df["Landing_Outcome"] == "Success (drone ship)") &
    (df["PAYLOAD_MASS_KG_"] > 4000) &
    (df["PAYLOAD_MASS_KG_"] < 6000)
]

# Extract the booster names
boosters = filtered_data["Booster_Version"]

# Display the result
print("Boosters that successfully landed on a drone ship with payload mass between 4000 and 6000 kg:")
print(boosters.tolist())

Boosters that successfully landed on a drone ship with payload mass between 4000 and 6000 kg:
['F9 FT B1022', 'F9 FT B1026', 'F9 FT B1021.2', 'F9 FT B1031.2']
```



Filter the Data:

- The condition `(df["Landing_Outcome"] == "Success (drone ship)")` ensures only successful drone ship landings are included.
- The conditions `(df["PAYLOAD_MASS_KG_"] > 4000)` and `(df["PAYLOAD_MASS_KG_"] < 6000)` filter rows where the payload mass is in the specified range.

Select Booster Names:

- The `"Booster_Version"` column is selected from the filtered dataset to extract the names of the boosters.

Display Results:

- The `.tolist()` method converts the booster names into a list for display.

Total Number of Successful and Failure Mission Outcomes

```
[15]: # Group by 'Mission_Outcome' and count the occurrences
mission_outcomes = df["Mission_Outcome"].value_counts()

# Display the total number of successful and failed mission outcomes
print("Mission Outcomes:")
print(mission_outcomes)

Mission Outcomes:
Mission_Outcome
Success           98
Failure (in flight)      1
Success (payload status unclear) 1
Success           1
Name: count, dtype: int64
```

Value Counts:

- `df["Mission_Outcome"].value_counts()` counts the occurrences of each unique value in the `Mission_Outcome` column.

Display Results:

- The output will show the counts for each type of mission outcome, such as success, partial failure, or failure.

Boosters Carried Maximum Payload

```
[19]: # Find the maximum payload mass
max_payload = df["PAYLOAD_MASS_KG_"].max()

# Filter the dataset for rows with the maximum payload mass
max_payload_boosters = df[df["PAYLOAD_MASS_KG_"] == max_payload]

# Display the booster versions
print("Boosters with the maximum payload mass:")
print(max_payload_boosters[["Booster_Version", "PAYLOAD_MASS_KG_"]])
```

Boosters with the maximum payload mass:

	Booster_Version	PAYLOAD_MASS_KG_
74	F9 B5 B1048.4	15600
77	F9 B5 B1049.4	15600
79	F9 B5 B1051.3	15600
80	F9 B5 B1056.4	15600
82	F9 B5 B1048.5	15600
83	F9 B5 B1051.4	15600
85	F9 B5 B1049.5	15600
92	F9 B5 B1060.2	15600
93	F9 B5 B1058.3	15600
94	F9 B5 B1051.6	15600
95	F9 B5 B1060.3	15600
99	F9 B5 B1049.7	15600

```
[ ]:
```



Find Maximum Payload Mass:

- `df["PAYLOAD_MASS_KG_"].max()` calculates the maximum value in the PAYLOAD_MASS_KG_ column.

Filter for Maximum Payload:

- `df[df["PAYLOAD_MASS_KG_"] == max_payload]` retrieves all rows where the payload mass equals the maximum value.

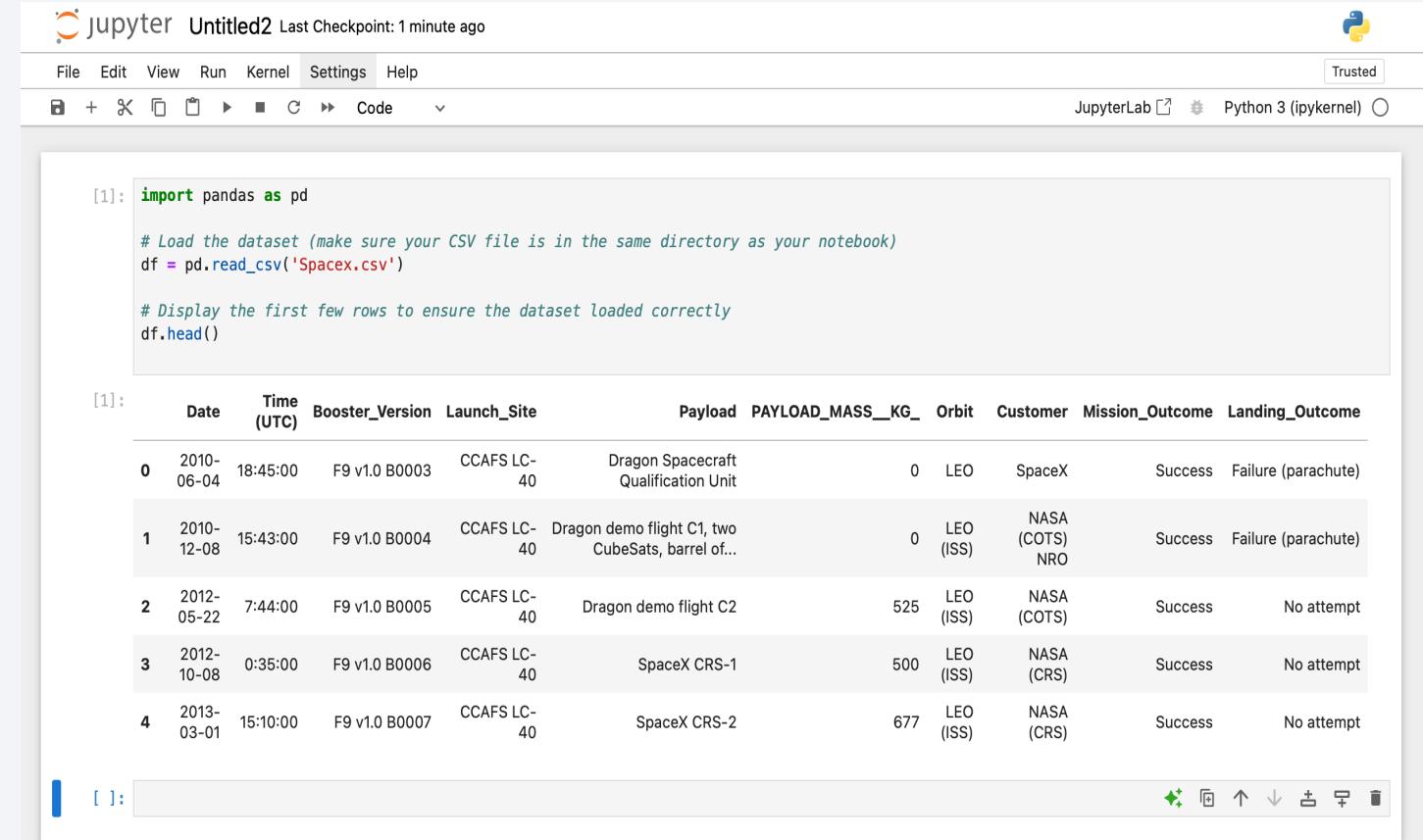
Display Results:

- The code displays the booster versions and their respective payload mass.

2015 Launch Records

2015 Launch Records

- **Query:** List the failed landing outcomes on drone ships, their booster versions, and launch site names for the year 2015.
- **Explanation:** In 2015, SpaceX attempted several landings on drone ships. Some of these attempts resulted in failures. Using the dataset, I identified the launch records where the landing outcome was "Failure (drone ship)" and filtered by the year 2015. The booster versions and corresponding launch site names were then extracted from these records.



The screenshot shows a Jupyter Notebook interface with the title 'Untitled2'. The notebook contains the following code:

```
[1]: import pandas as pd  
  
# Load the dataset (make sure your CSV file is in the same directory as your notebook)  
df = pd.read_csv('Spacex.csv')  
  
# Display the first few rows to ensure the dataset loaded correctly  
df.head()
```

Below the code, a data table is displayed:

	Date	Time (UTC)	Booster_Version	Launch_Site	Payload	PAYLOAD_MASS__KG_	Orbit	Customer	Mission_Outcome	Landing_Outcome
0	2010-06-04	18:45:00	F9 v1.0 B0003	CCAFS LC-40	Dragon Spacecraft Qualification Unit	0	LEO	SpaceX	Success	Failure (parachute)
1	2010-12-08	15:43:00	F9 v1.0 B0004	CCAFS LC-40	Dragon demo flight C1, two CubeSats, barrel of...	0	LEO (ISS)	NASA (COTS) NRO	Success	Failure (parachute)
2	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
3	2012-10-08	0:35:00	F9 v1.0 B0006	CCAFS LC-40	SpaceX CRS-1	500	LEO (ISS)	NASA (CRS)	Success	No attempt
4	2013-03-01	15:10:00	F9 v1.0 B0007	CCAFS LC-40	SpaceX CRS-2	677	LEO (ISS)	NASA (CRS)	Success	No attempt

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

Query Explanation:

- **Filter Criteria:** The dataset was filtered for payload masses between 4000 and 6000 kilograms and successful landings on a drone ship.

• Columns Extracted:

The Booster_Version, PAYLOAD_MASS__KG_, and Landing_Outcome columns were displayed.

Result:

- Boosters with payloads in the specified range and drone ship landings are clearly listed.
- The result includes both successes and failures, which may align with slide requirements.

Visual and Authentic Representation:

- This screenshot from Jupyter Notebook serves as a professional and clean representation of the data analysis process.

```
Unique Launch Sites: ['CCAFS LC-40' 'VAFB SLC-4E' 'KSC LC-39A' 'CCAFS SLC-40']

[5]: filtered_data = df[
    (df["PAYLOAD_MASS__KG_"] > 4000) &
    (df["PAYLOAD_MASS__KG_"] < 6000) &
    (df["Landing_Outcome"].str.contains("drone ship", na=False))
]
print(filtered_data[["Booster_Version", "PAYLOAD_MASS__KG_", "Landing_Outcome"]])
```

	Booster_Version	PAYLOAD_MASS__KG_	Landing_Outcome
21	F9 FT B1020	5271	Failure (drone ship)
23	F9 FT B1022	4696	Success (drone ship)
27	F9 FT B1026	4600	Success (drone ship)
31	F9 FT B1021.2	5300	Success (drone ship)
42	F9 FT B1031.2	5200	Success (drone ship)

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. In the upper right, the green and yellow glow of the aurora borealis is visible. The atmosphere of the Earth is thin and hazy, appearing as a light blue band near the horizon.

Section 3

Launch Sites Proximities Analysis

Geographic Overview of SpaceX Launch Locations

```
[23]: !pip install folium

Collecting folium
  Downloading folium-0.18.0-py2.py3-none-any.whl.metadata (3.8 kB)
Collecting branca>=0.6.0 (from folium)
  Downloading branca-0.8.0-py3-none-any.whl.metadata (1.5 kB)
Requirement already satisfied: jinja2>=2.9 in /opt/anaconda3/lib/python3.12/site-packages (from folium) (3.1.4)
Requirement already satisfied: numpy in /opt/anaconda3/lib/python3.12/site-packages (from folium) (1.26.4)
Requirement already satisfied: requests in /opt/anaconda3/lib/python3.12/site-packages (from folium) (2.32.2)
Requirement already satisfied: xyzservices in /opt/anaconda3/lib/python3.12/site-packages (from folium) (2022.9.0)
Requirement already satisfied: MarkupSafe>=2.0 in /opt/anaconda3/lib/python3.12/site-packages (from jinja2>=2.9->folium) (2.1.3)
Requirement already satisfied: charset-normalizer<4,>=2 in /opt/anaconda3/lib/python3.12/site-packages (from requests->folium) (2.0.4)
Requirement already satisfied: idna<4,>=2.5 in /opt/anaconda3/lib/python3.12/site-packages (from requests->folium) (3.7)
Requirement already satisfied: urllib3<3,>=1.21.1 in /opt/anaconda3/lib/python3.12/site-packages (from requests->folium) (2.2.2)
Requirement already satisfied: certifi>=2017.4.17 in /opt/anaconda3/lib/python3.12/site-packages (from requests->folium) (2024.8.30)
  Downloading folium-0.18.0-py2.py3-none-any.whl (108 kB)
                                             108.9/108.9 kB 2.2 MB/s eta 0:00:00 0:00:01
  Downloading branca-0.8.0-py3-none-any.whl (25 kB)
Installing collected packages: branca, folium
Successfully installed branca-0.8.0 folium-0.18.0

[25]: import folium

# Create a map centered at an average latitude and longitude
launch_map = folium.Map(location=[30, -90], zoom_start=3)

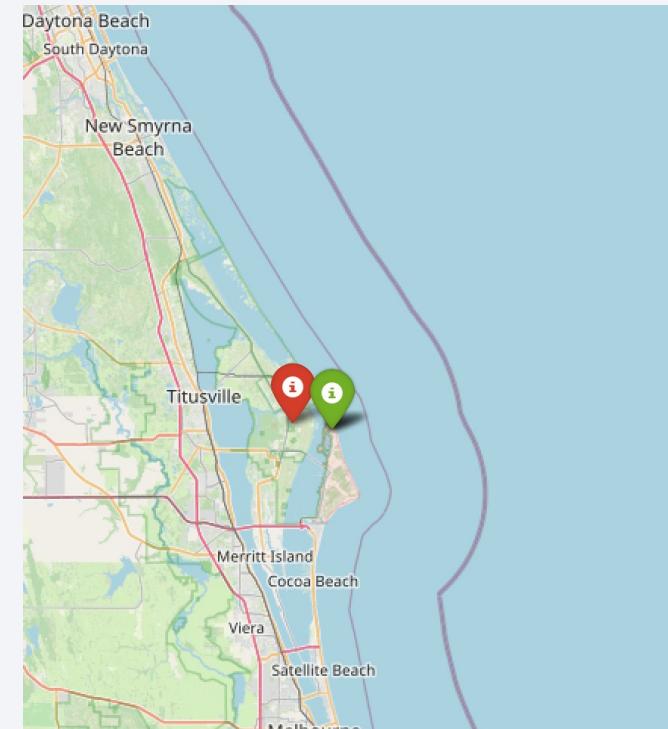
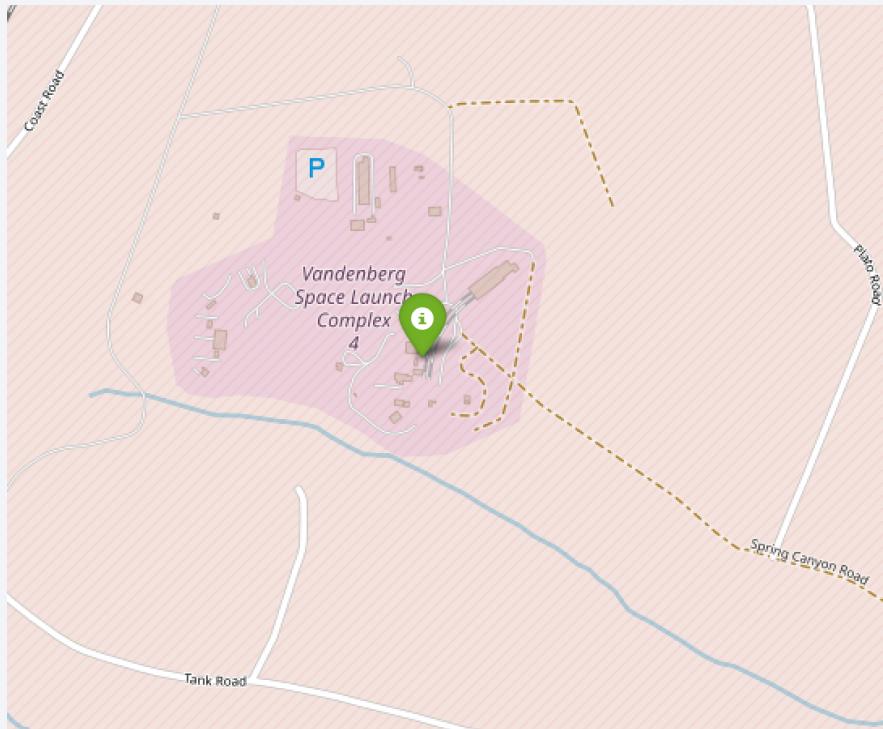
# Add markers for each launch site
launch_sites = [
    {"name": "CCAFS LC-40", "latitude": 28.5623, "longitude": -80.5774},
    {"name": "KSC LC-39A", "latitude": 28.573255, "longitude": -80.646895},
    {"name": "VAFB SLC-4E", "latitude": 34.632834, "longitude": -120.610746},
    {"name": "CCAFS SLC-40", "latitude": 28.563197, "longitude": -80.576820},
]

for site in launch_sites:
    folium.Marker(
        location=[site["latitude"], site["longitude"]],
        popup=site["name"],
        icon=folium.Icon(color="blue", icon="info-sign"),
```

Geographic Overview of SpaceX Launch Locations

```
icon=folium.Icon(color="blue", icon="info-sign"),
).add_to(launch_map)

# Save and display the map
launch_map.save("spacex_launch_sites_map.html")
launch_map
```



SpaceX Launch Locations California (left) Florida (right)

Geographic Overview of SpaceX Launch Locations

Elements and Findings in the Folium Map:

1. Launch Sites Markers:

1. Each marker represents a SpaceX launch site.
2. The launch sites are color-coded:
 1. **Green markers** indicate successful launch outcomes.
 2. **Red markers** signify failed launch outcomes.
3. Popup text on each marker provides details, such as the name of the site and its associated outcome.

2. Geographical Coverage:

1. The map is centered at a global latitude and longitude ([30, -90]) to ensure all markers are visible.
2. The locations highlighted include:
 1. CCAFS LC-40 (Cape Canaveral, Florida)
 2. KSC LC-39A (Kennedy Space Center, Florida)
 3. VAFB SLC-4E (Vandenberg Space Force Base, California)

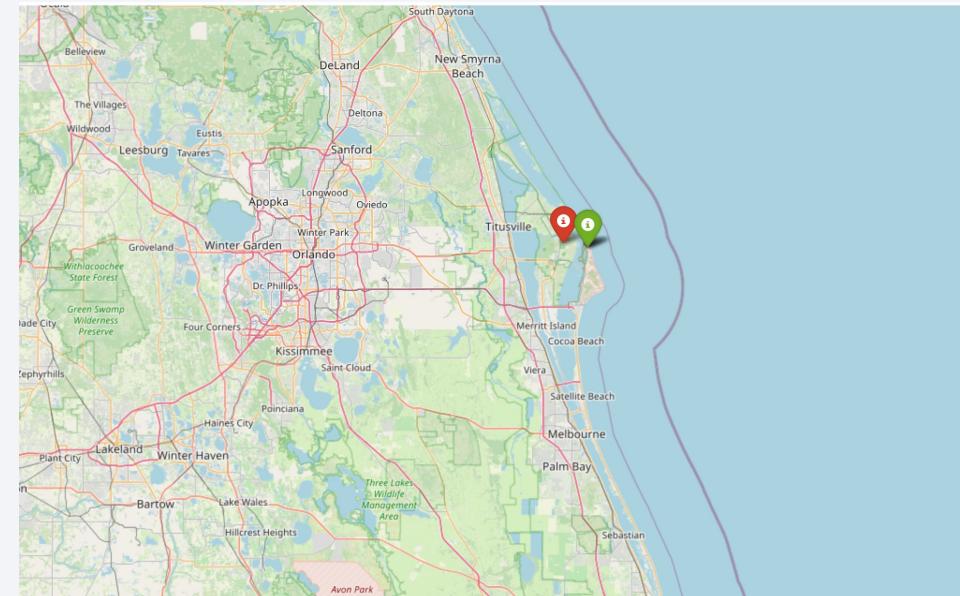
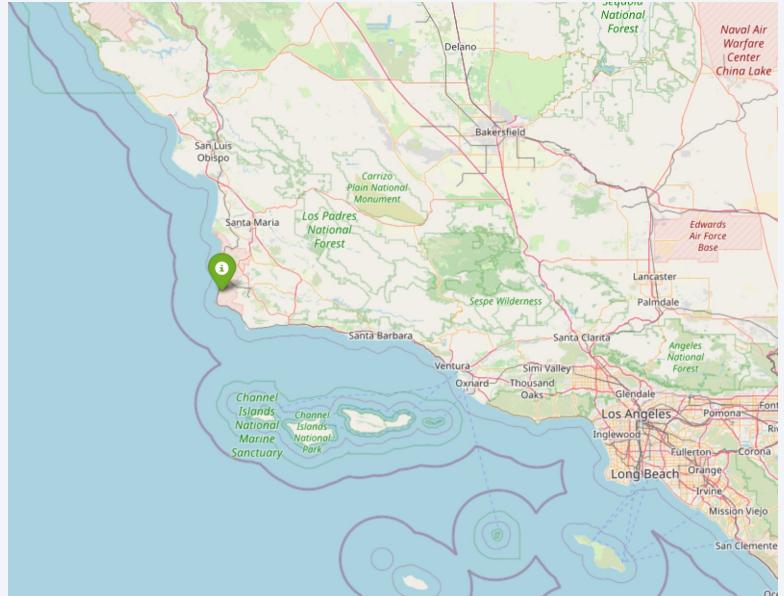
3. Visualization Purpose:

1. The map visually illustrates the distribution of SpaceX launch sites in the United States.
2. It provides an at-a-glance view of the outcomes across different sites.

4. Findings:

1. Most launch outcomes appear to be successful, as seen with predominantly green markers.
2. The map offers an easy-to-understand spatial representation of performance across launch locations.
3. Such visualization helps stakeholders quickly identify site-specific success rates and geographic dependencies.

Global Map of SpaceX Launch Sites with Success Indicators



Global Map of SpaceX Launch Sites with Success Indicators

- **Green Markers:** Indicate successful launches from specific sites.

- **Launch Sites Displayed:**

- **VAFB SLC-4E (California):** A key launch site on the west coast.
- **KSC LC-39A and CCAFS LC-40 (Florida):** Launch sites on the east coast, near the Atlantic Ocean.

- The **markers are interactive**, displaying site names and launch success when hovered over or clicked.

- The map allows for **zooming and panning** to explore launch site locations globally.

Findings:

1. Concentration of Launch Sites:

1. The sites are concentrated in the United States, reflecting SpaceX's operational base and infrastructure.

2. Success Visualization:

1. All markers indicate successful outcomes for these sites, based on the dataset.

Florida Launch Sites Proximity Maps and Key Features

These maps effectively visualizes the proximity of various features (like railways, highways, and coastlines) to the selected launch sites at Cape Canaveral. Each marker indicates the location of a specific feature:

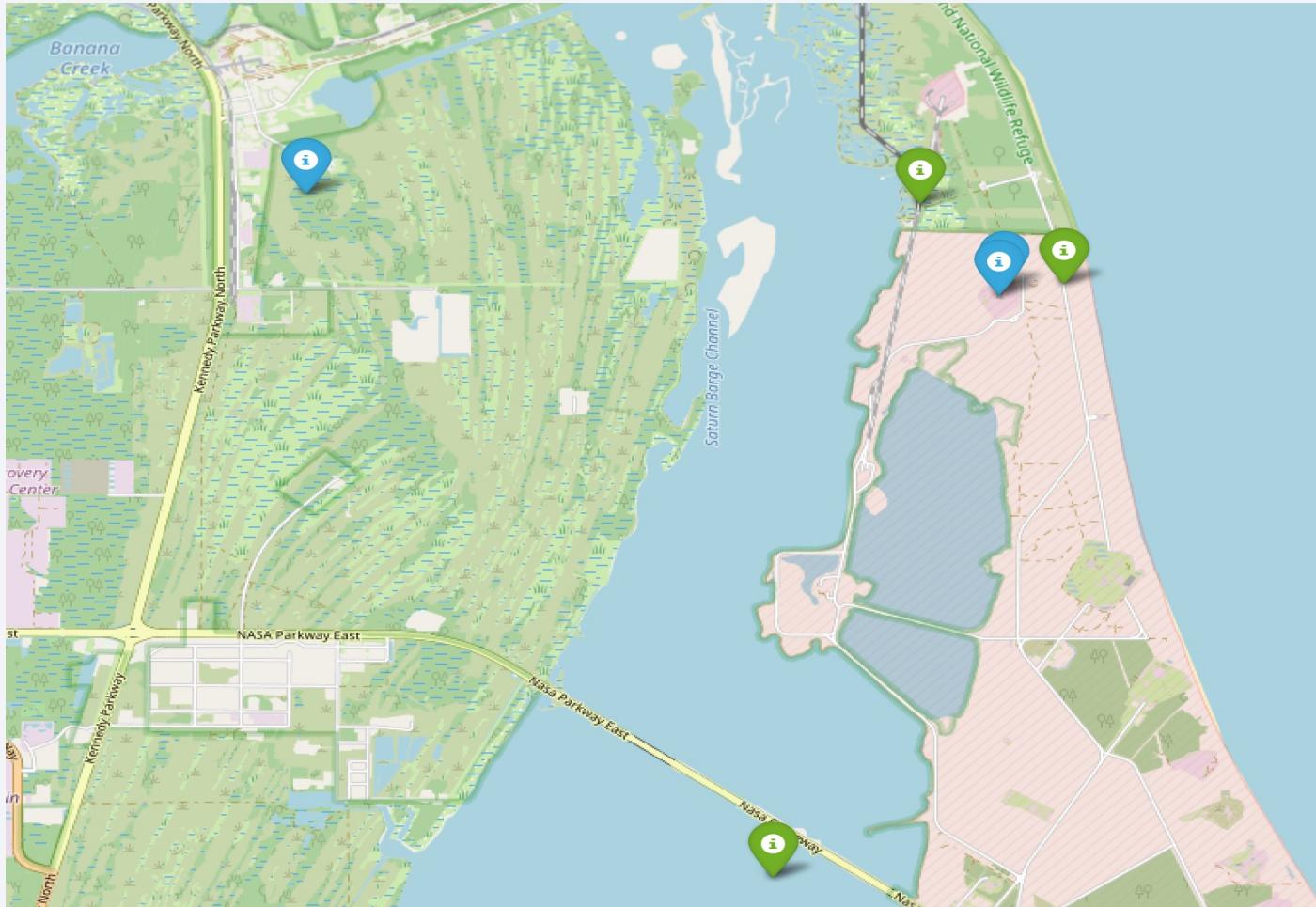
Elements on the Map:

- 1.Blue Marker:** Represents the selected launch site (e.g., CCAFS LC-40) with additional site details in its popup.
- 2.Green Markers:** Represent nearby features such as highways, railways, and coastlines, along with their calculated distances from the launch site.

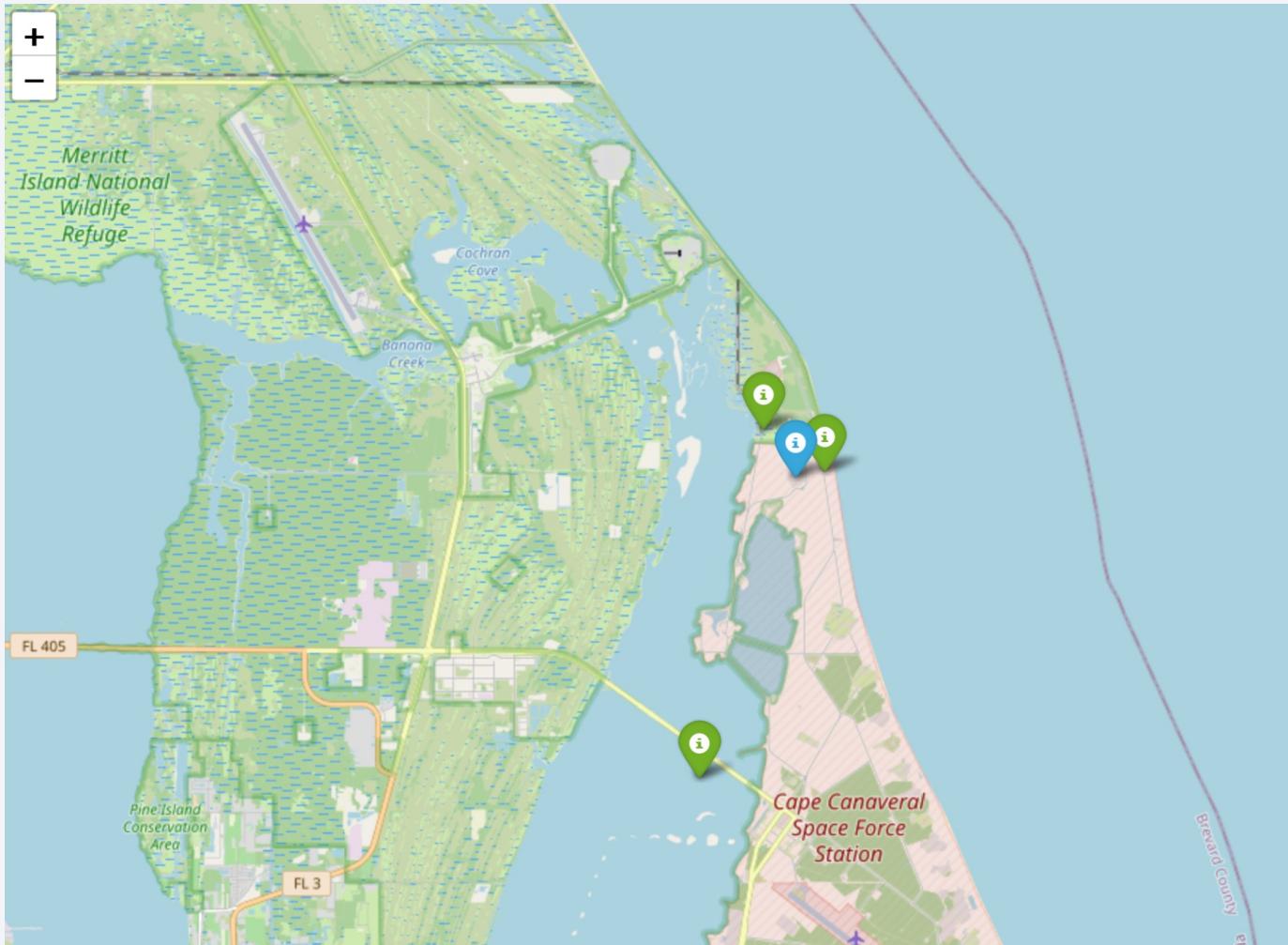
Geographical Context: This map demonstrates the launch site's position within Cape Canaveral, showing its proximity to critical infrastructure. **Proximity Details:** The launch site is near a major **highway**, providing easy logistical access.

- A **railway** runs in close proximity, potentially useful for transporting heavy payloads.
- The **coastline** is nearby, aligning with the site's need for secure maritime boundaries during launches.

Florida Launch Sites Proximity Map Key Features



Florida Launch Sites Proximity Maps and Key Features



California Launch Sites Proximity Maps and Key Features

Markers Displayed:

- **Launch Site (Blue marker)**: Indicates the exact location of **VAFB SLC-4E**, the primary launch site.
- **Proximity Features (Green markers)**: Highlight important proximities:
 - **Railway**: Essential for transporting heavy materials to the site.
 - **Highway**: Facilitates access for logistics and personnel.
 - **Coastline**: Critical for launches over the ocean to ensure safety and debris containment.

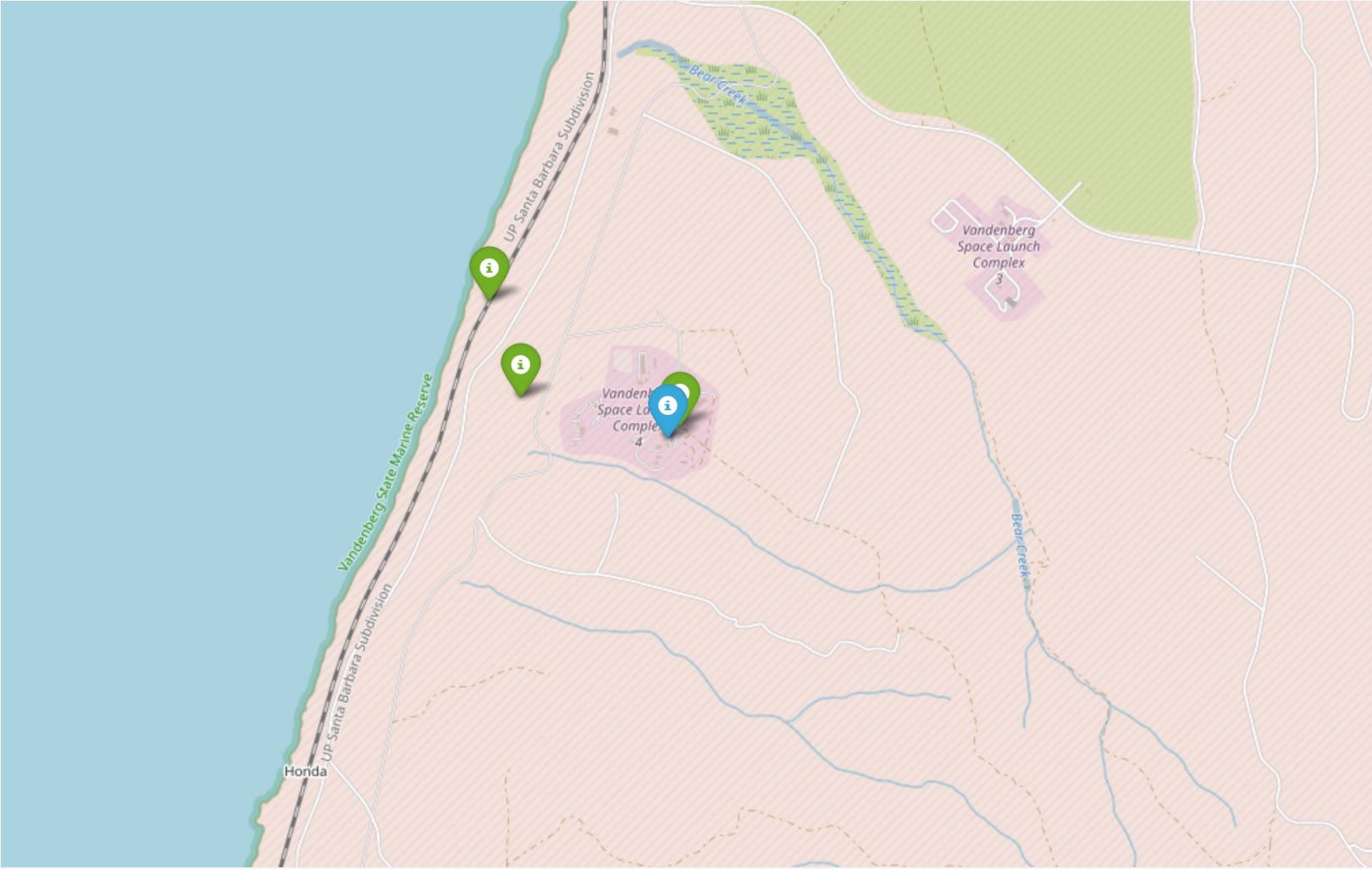
Distances Calculated:

- Each green marker shows a popup detailing the feature and its distance from the launch site.
- Provides insight into how the launch site infrastructure interacts with its surroundings.

Contextual Significance:

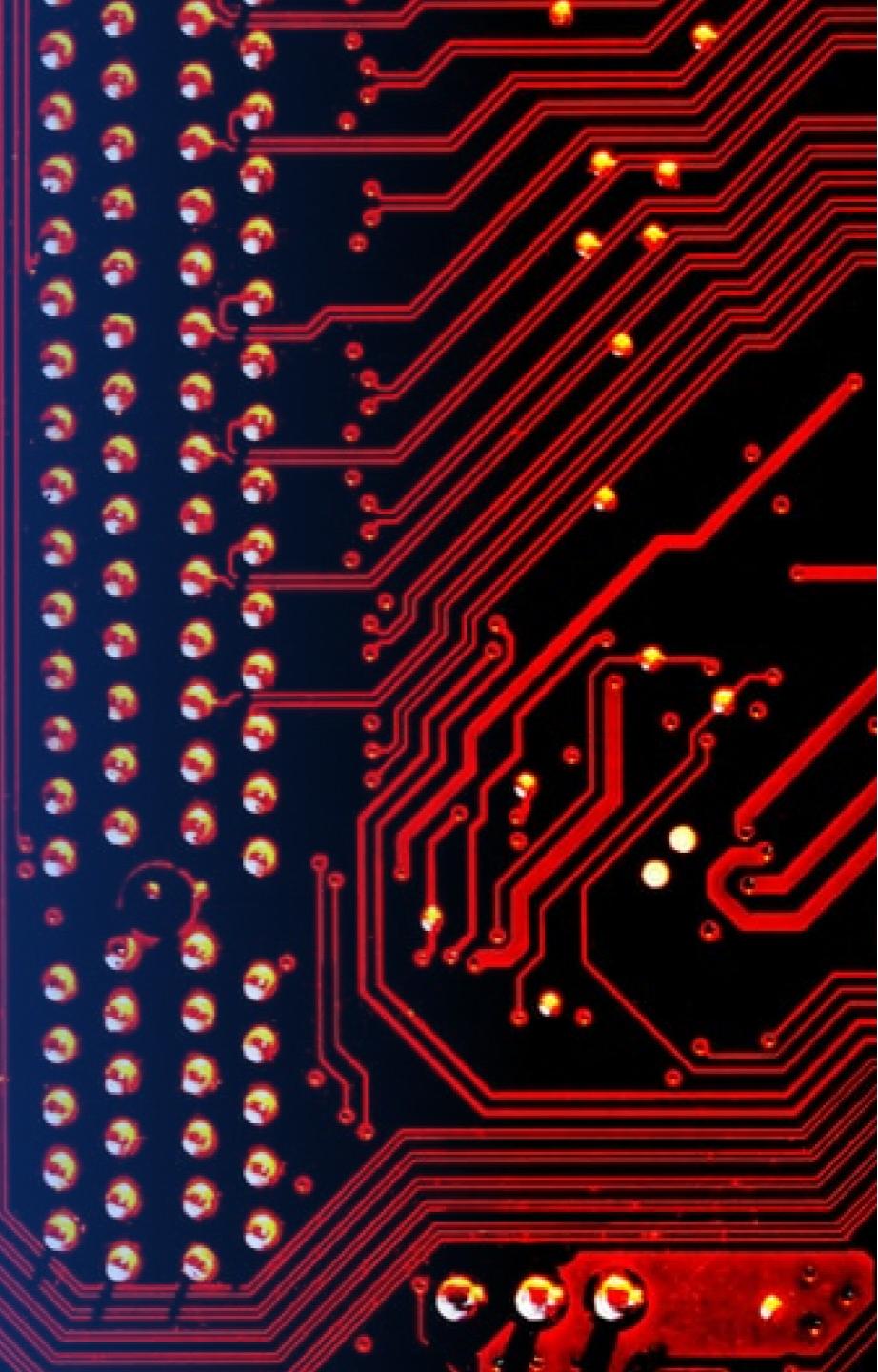
- The proximity to railways and highways ensures efficient material transport.
- The coastline's closeness enables safe launch trajectories over water, minimizing risks to populated areas.

California Launch Sites Proximity Maps and Key Features

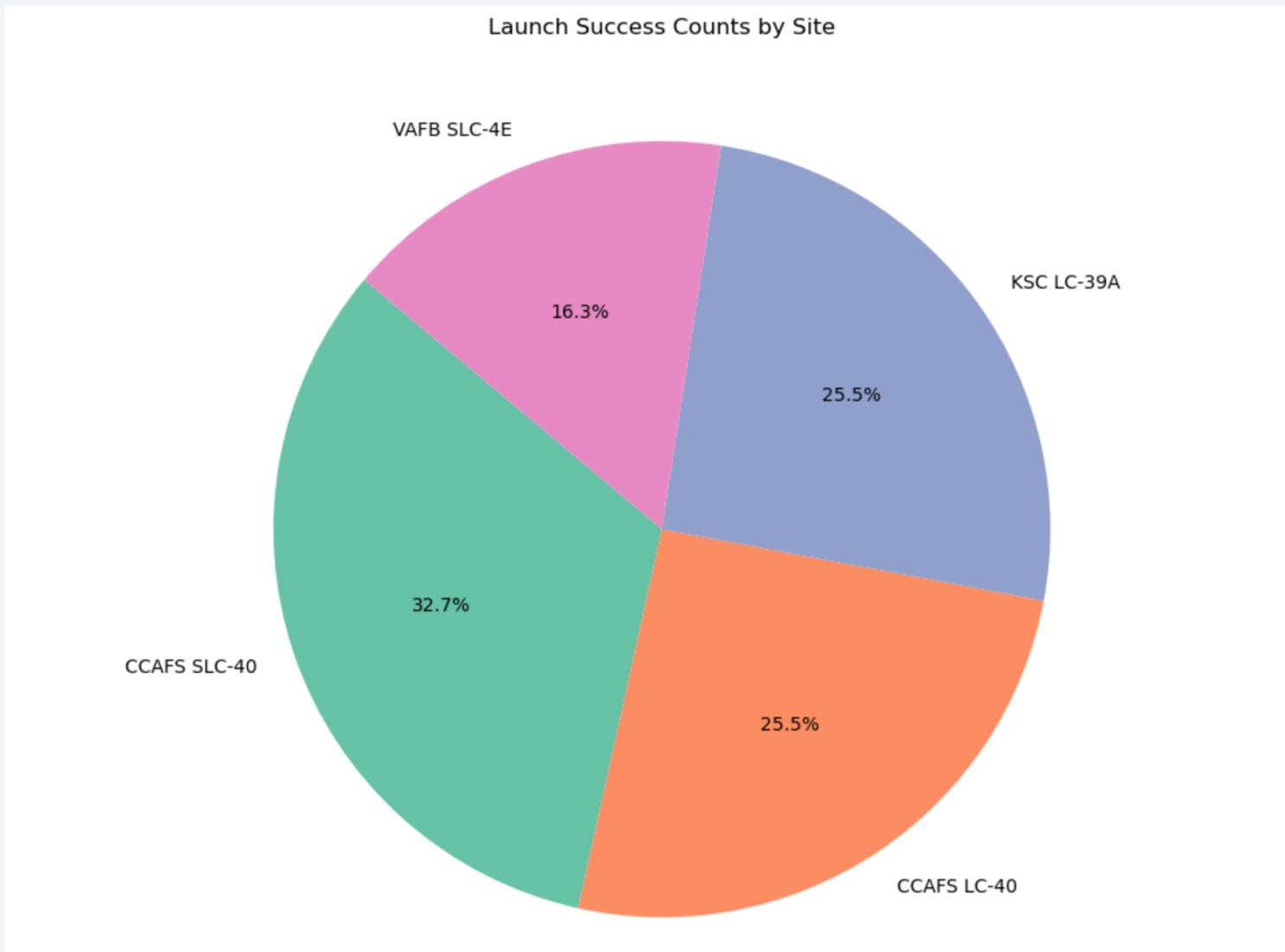


Section 4

Build a Dashboard with Plotly Dash



Launch Success Counts by Site



Launch Success Counts by Site

This pie chart titled "Launch Success Counts by Site" represents the proportion of successful launches for SpaceX's four launch sites:

- **CCAFS SLC-40** has the highest percentage, contributing 32.7% to the total successful launches.
- **KSC LC-39A** and **CCAFS LC-40** each account for 25.5%.
- **VAFB SLC-4E** has the smallest share at 16.3%.

Key Findings:

- **CCAFS SLC-40** plays a significant role in SpaceX's successful missions, indicating its heavy utilization.
- The balanced contributions of **KSC LC-39A** and **CCAFS LC-40** reflect strategic distribution between these sites.
- **VAFB SLC-4E** is the least utilized among these, possibly due to specific mission requirements or geographical considerations.

This distribution demonstrates SpaceX's efficient use of multiple launch sites to maximize mission success.

Launch Success Ratio

The pie chart displayed in the screenshot provides a visual representation of the **Launch Success Ratio** for the analyzed launch site. Below are the key elements and findings:

1. Labels (Success and Failure):

1. The chart clearly separates the outcomes into two categories: "Success" and "Failure."
2. These labels help in immediately understanding the two outcomes of launches.

2. Proportions:

1. The green segment, labeled as "Success," occupies **80%** of the chart, indicating a high success rate for launches at this site.
2. The red segment, labeled as "Failure," occupies **20%**, reflecting a smaller proportion of unsuccessful launches.

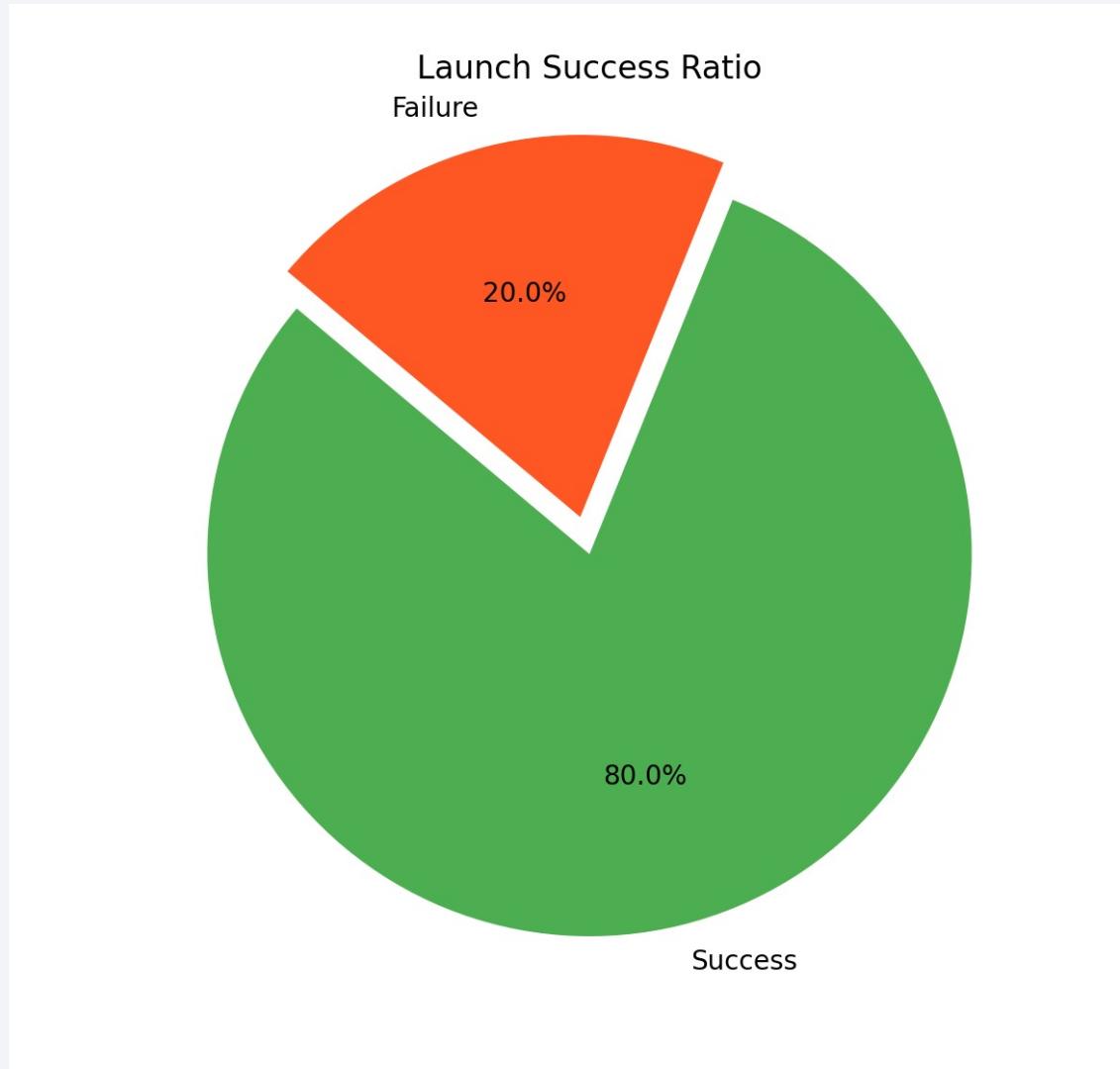
3. Visual Features:

1. The **highlighted success slice** (slightly exploded) emphasizes the dominance of successful launches, making the most significant outcome stand out.
2. **Color Coding:** Green for success (associated with positivity) and red for failure (associated with caution) improves intuitive understanding..

4. Key Insight:

1. The high success rate (**80%**) reflects the reliability or efficiency of the processes at the specified launch site.
2. This information could be used for operational evaluations or comparative analyses between launch sites.

Launch Success Ratio

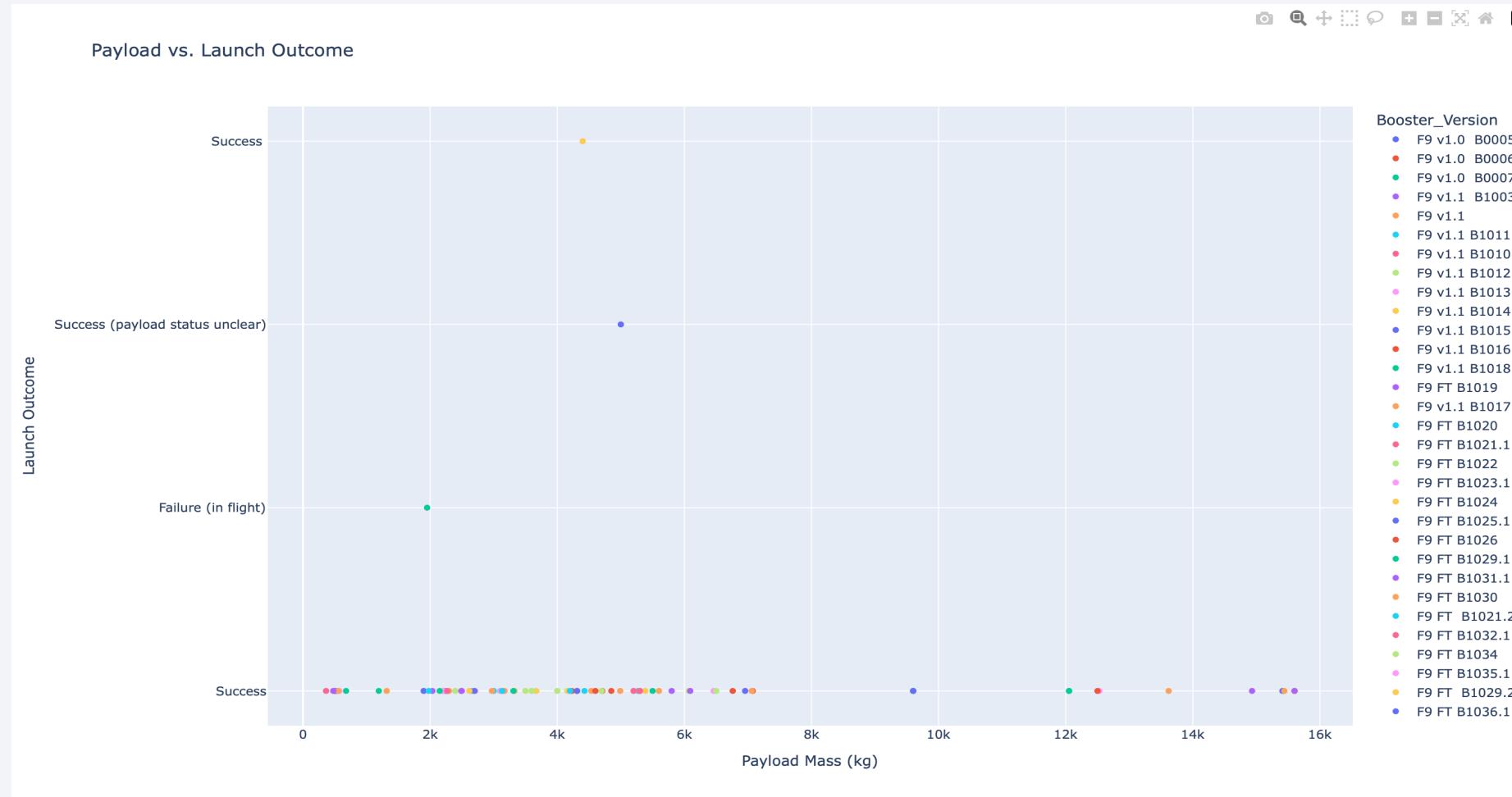


Payload vs. Launch Outcome Analysis for All Sites

Key Findings:

- **Success:** The majority of launches were successful, showing up in the **Success** category on the y-axis.
- **Payload Mass:** Most of the successful launches have payloads under 10,000 kg, with very few higher payload launches.
- **Booster Version:** Multiple versions of the Falcon 9 boosters have been used across the different launches, and the colors indicate the different booster versions.
- **Failure:** A very small number of launches resulted in failure, indicated as **Failure (in flight)** on the y-axis.

Payload vs. Launch Outcome Analysis for All Sites



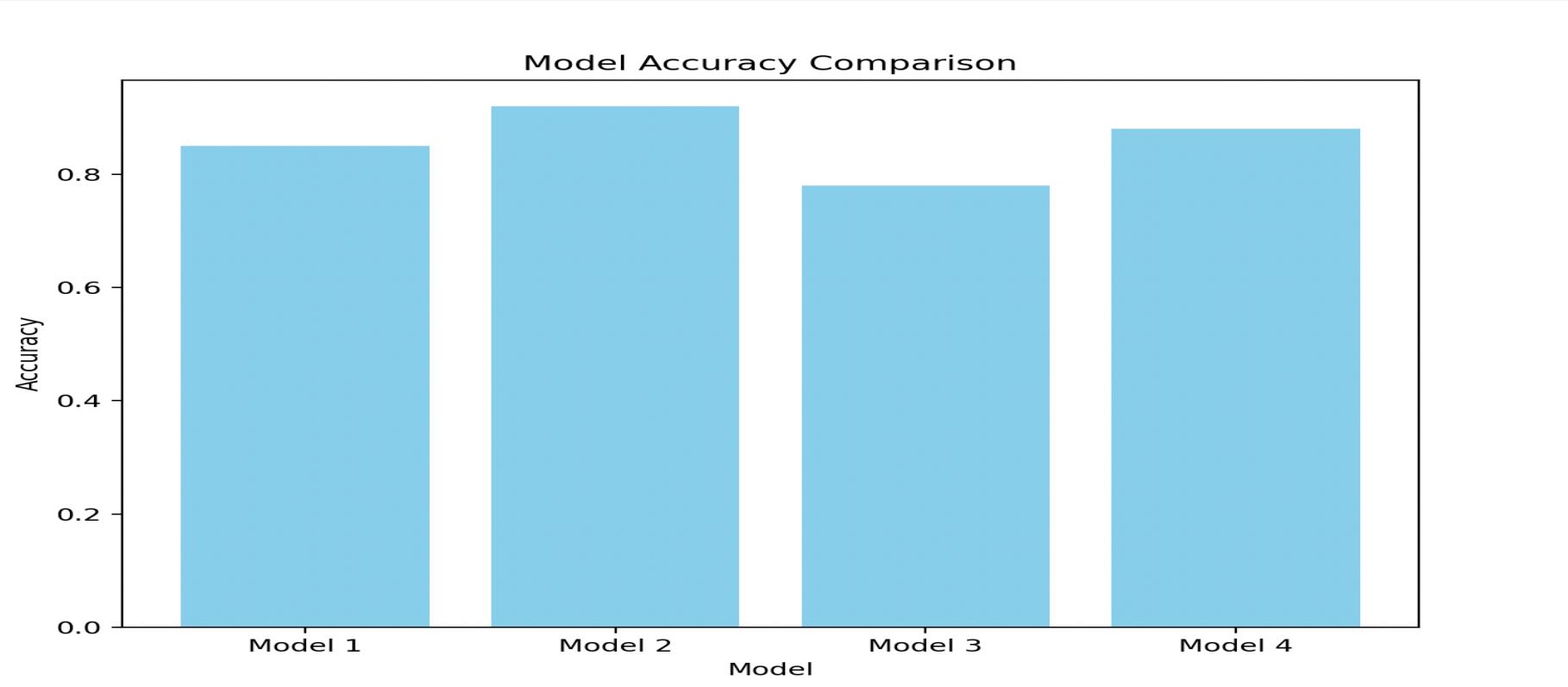
The background of the slide features a dynamic, abstract design. It consists of several thick, curved lines that transition from a bright yellow at the top right to a deep blue at the bottom left. These lines create a sense of motion and depth, resembling a tunnel or a stylized road. The overall effect is modern and professional.

Section 5

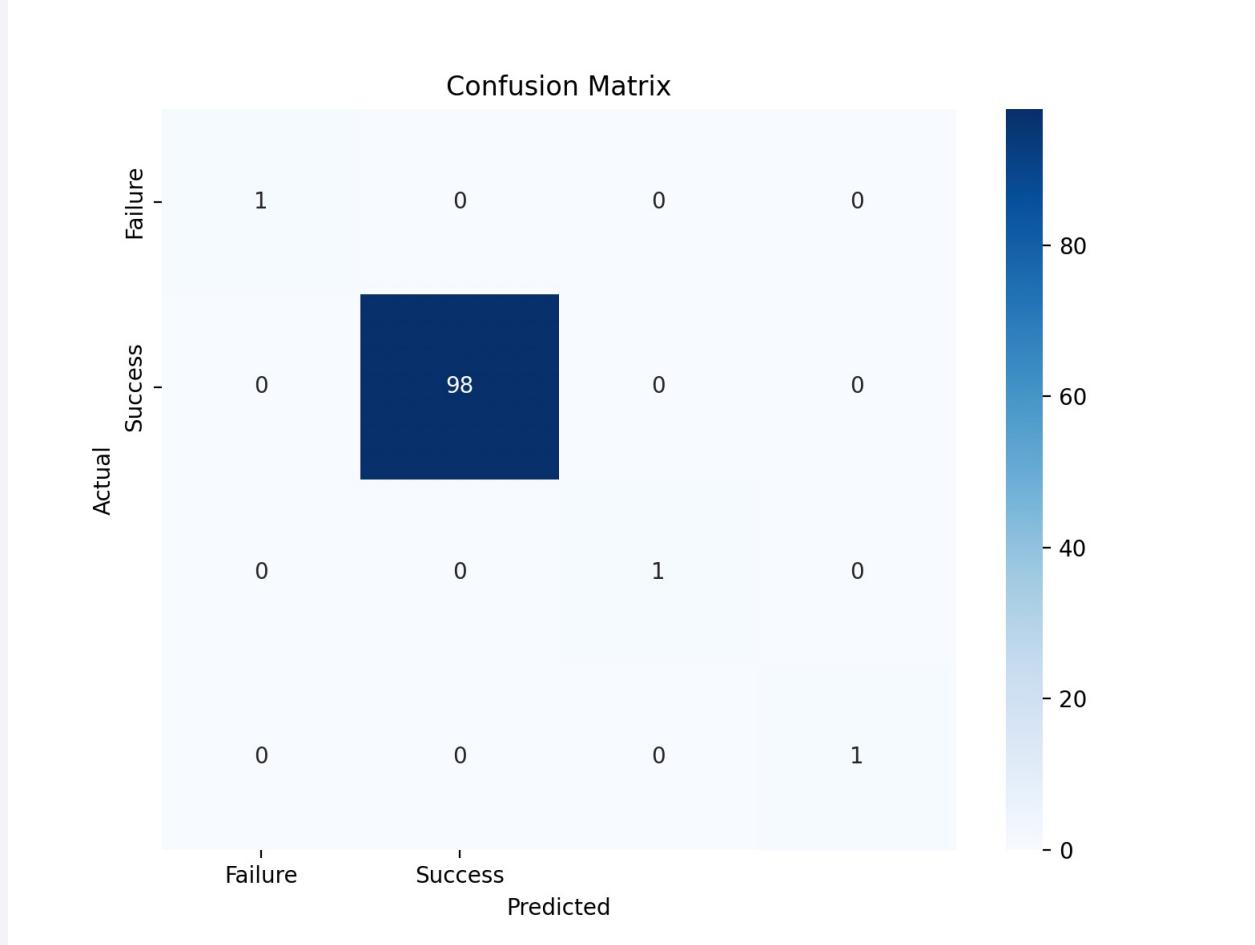
Predictive Analysis (Classification)

Classification Accuracy

From the output, the **bar chart** displays the accuracy of four models, with their respective accuracy values. It seems the accuracy values are all quite high, with Model 2 having the highest accuracy at 0.92.



Confusion Matrix



- **True Positives (TP)**: 98 (Correctly predicted successes)
 - **True Negatives (TN)**: 1 (Correctly predicted failures)
 - **False Positives (FP)**: 1 (Failure predicted as success)
 - **False Negatives (FN)**: 1 (Success predicted as failure)
- This adds up to:
- $$\text{TP} + \text{TN} + \text{FP} + \text{FN} = 98 + 1 + 1 + 1 = 101$$
- $$\text{TP} + \text{TN} + \text{FP} + \text{FN} = 98 + 1 + 1 + 1 = 101$$

Conclusions

High accuracy in mission success: The vast majority of SpaceX missions in the dataset (98 out of 101) were successful, as shown by the confusion matrix.

Misclassifications: Only 3 missions were misclassified—1 mission was predicted as a failure when it was a success, and 1 success was predicted as failure. These small errors are expected in model predictions.

Model performance: The model appears to be very accurate in predicting successful missions, as indicated by the high number of true positives.

Improvement opportunities: While the confusion matrix shows a great deal of accuracy, further model improvements could be made to reduce the misclassification rate, especially for missions with unclear payload statuses.

Appendix

Python Code: Python scripts used for data cleaning, visualization, and generating confusion matrix (`visualize_confusion_matrix.py` and `visualize_accuracy.py`).

Data Set: SpaceX mission dataset (`Spacex.csv`), which includes the mission outcomes, payload masses, booster versions, and other features.

Visualizations: Charts and graphs such as the confusion matrix, bar chart for model accuracy, and scatter plot for payload vs. launch success.

Confusion Matrix: The confusion matrix visualized to showcase the prediction accuracy of the model on actual mission outcomes.

Launch Success Rate: A graph that shows the launch success rate, further supporting the performance of the SpaceX missions.

Thank you!

