# Outline

- Executive Summary

- Introduction

- Methodology

- Results

- Conclusion

- Appendix

# Executive Summary

- The following methodologies were used to analyze data:

  - Data Collection using web scraping and SpaceX API

  - Exploratory Data Analysis (EDA), including data wrangling, data visualization and interactive visual analytics

  - Machine Learning Prediction.

- Summary of all results:

  - It was possible to collect valuable data from public sources

  - EDA allowed to identify which features are the best to predict success of launchings

  - Machine Learning Prediction showed the best model to predict which characteristics are important to drive this opportunity by the best way, using all collected data.

# Introduction

- Project background and context: This project analyzes historical SpaceX launch data to identify key factors that influence mission success. By applying exploratory data analysis and machine learning, we aim to predict launch outcomes and support decision-making in the aerospace industry. It serves as the final assignment for the Coursera Applied Data Science Capstone.

- Problems you want to find answers:

  - What features impact launch success the most?

  - Can we accurately predict launch outcomes?

  - Which machine learning model performs best?

  - How do launch sites differ in success rates?

Section 1

# Methodology

# Methodology

Executive Summary

- Data collection methodology:

  - Describe how data was collected

- Perform data wrangling

  - Describe how data was processed

- Perform exploratory data analysis (EDA) using visualization and SQL

- Perform interactive visual analytics using Folium and Plotly Dash

- Perform predictive analysis using classification models

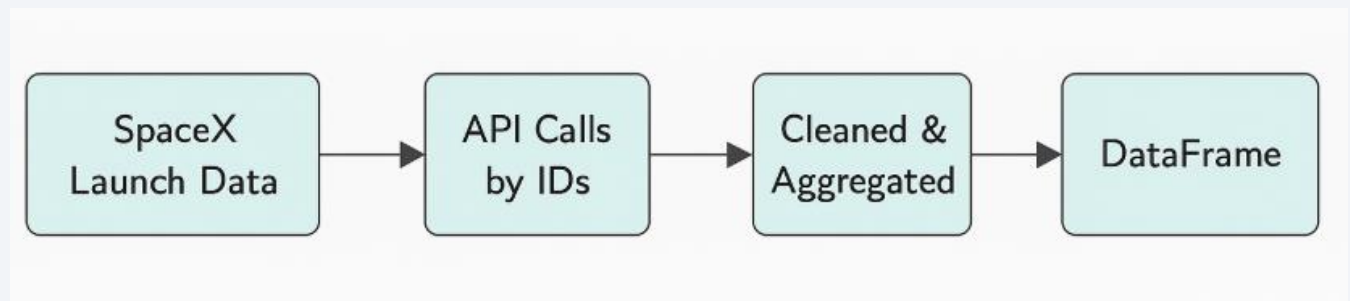  - How to build, tune, evaluate classification models

# Data Collection

- Data was collected using the SpaceX public REST API (https://api.spacexdata.com/v4).

- Python requests was used to query launch, rocket, payload, and launchpad data.

- Multiple helper functions were implemented to extract:

  - BoosterVersion from rocket IDs

  - LaunchSite, Longitude, and Latitude from launchpad IDs

  - PayloadMass and Orbit from payload IDs

  - Block, ReusedCount, Outcome, etc. from core details.

# Data Collection – SpaceX API

- Present your data collection with SpaceX REST calls using key phrases and flowcharts

- Add the GitHub URL of the completed SpaceX API calls notebook (must include completed code cell and outcome cell), as an external reference and peer-review purpose:
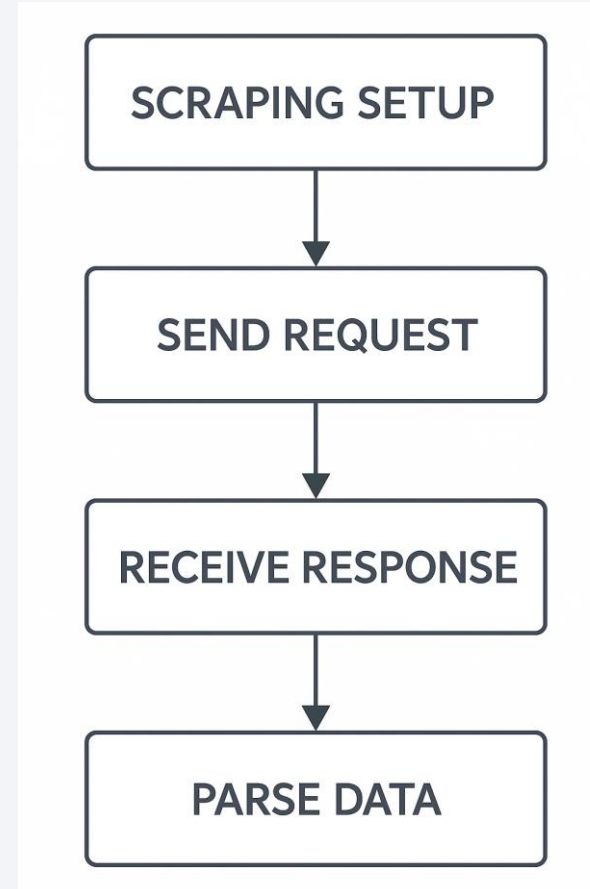
https://github.com/tranhoangthienphu/Applied-Data-Science-Capstone/blob/cf92e8c938f8d0ced1b6073477af64fbad6e9471/jupyter-labs-spacex-data-collection-api.ipynb

SpaceX Launch Data → API Calls by IDs → Cleaned & Aggregated → DataFrame

# Data Collection - Scraping

- Present your web scraping process using key phrases and flowcharts

- Add the GitHub URL of the completed web scraping notebook, as an external reference and peer-review purpose

https://github.com/tranhoangthienphu/Applied-Data-Science-Capstone/blob/cf92e8c938f8d0ced1b6073477af64fbad6e9471/jupyter-labs-webscraping.ipynb



SCRAPING SETUP

SEND REQUEST

RECEIVE RESPONSE

PARSE DATA

# Data Wrangling

- Describe how data were processed

- You need to present your data wrangling process using key phrases and flowcharts

- Add the GitHub URL of your completed data wrangling related notebooks, as an external reference and peer-review purpose

# EDA with Data Visualization

- Summarize what charts were plotted and why you used those charts:

  - Bar charts were used to display the number of launches by orbit type and launch site, helping identify the most frequently used orbits and locations.

  - Pie charts showed the proportion of successful vs. failed launches, offering a quick understanding of SpaceX's overall success rate.

  - Scatter plots illustrated the relationship between payload mass and launch success, helping evaluate if heavier payloads affect outcomes.

  - Histograms were used to show the distribution of payload mass across different missions, highlighting common payload ranges.

- Add the GitHub URL of your completed EDA with data visualization notebook, as an external reference and peer-review purpose:

https://github.com/tranhoangthienphu/Applied-Data-Science-Capstone/blob/cf92e8c938f8d0ced1b6073477af64fbad6e9471/EDA%20with%20Visualization%20Lab.ipynb

# EDA with SQL

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

  - Retrieved average payload mass carried by F9 v1.1 boosters to analyze performance.

  - Identified the date of the first successful ground pad landing.

  - Listed boosters that succeeded in drone ship landings with specific payload conditions.

  - Counted total successful and failed mission outcomes for all launches.

  - Extracted boosters with the maximum payload using subqueries and aggregation.

  - Filtered 2015 launch failures on drone ships by month and site.

  - Ranked landing outcomes by frequency between specific dates using SQL sorting and grouping.

- Add the GitHub URL of your completed EDA with SQL notebook, as an external reference and peer-review purpose:
  https://github.com/tranhoangthienphu/SQL.ipynb

# Build an Interactive Map with Folium

- Added markers to display all SpaceX launch sites, and circles to show the proximity around each site.

- Used polylines to visually connect launch sites with nearest coastlines.

- Purpose: to provide spatial insights on launch distribution, proximity to ocean, and success probability by location.

- GitHub URL: https://github.com/tranhoangthienphu/Applied-Data-Science-Capstone/blob/cf92e8c938f8d0ced1b6073477af64fbad6e9471/Interactive%20Visual%20Analytics%20with%20Folium%20lab_launch_site_location.ipynb

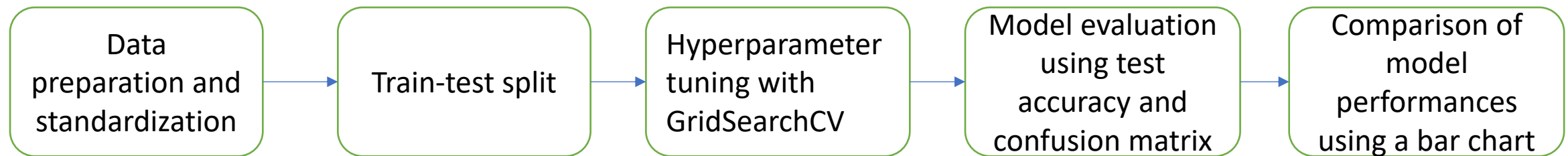# Build a Dashboard with Plotly Dash

- **Dashboard Features:**

    - Dropdown menu to select launch site (all sites or specific).

    - Interactive pie chart to show success rates by site or by success/failure for selected site.

    - Range slider to filter payload mass range.

    - Interactive scatter plot to explore correlation between payload and launch outcome, color-coded by booster version.

- **Why These Were Used:**

    - Enable users to explore data dynamically.

    - Support better understanding of launch performance by site and payload.

    - Visualize relationships and trends for improved decision-making.

- **GitHub URL:** https://github.com/tranhoangthienphu/Applied-Data-Science-Capstone/blob/cf92e8c938f8d0ced1b6073477af64fbad6e9471/spacex-dash-app.py

# Predictive Analysis (Classification)

- Summarize how you built, evaluated, improved, and found the best performing classification model:

  - Data was preprocessed and standardized.

  - Several classification models were trained and tuned using GridSearchCV: Logistic Regression

  - Support Vector Machine (SVM), Decision Tree, and K-Nearest Neighbors (KNN).

  - Each model's hyperparameters were optimized with 10-fold cross-validation.

  - Model performance was evaluated using accuracy scores and confusion matrices.

# Predictive Analysis (Classification)

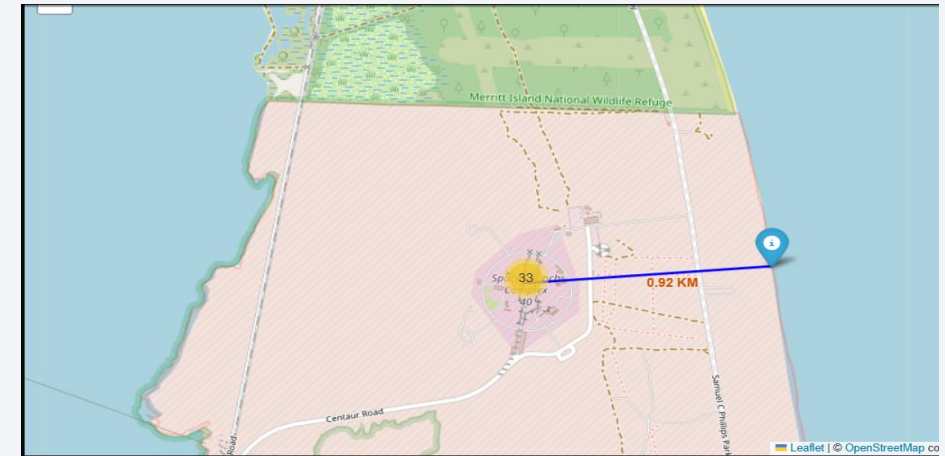- You need present your model development process using key phrases and flowchart:

| Data preparation and standardization | → | Train-test split | → | Hyperparameter tuning with GridSearchCV | → | Model evaluation using test accuracy and confusion matrix | → | Comparison of model performances using a bar chart |
|---|---|---|---|---|---|---|---|---|

- Add the GitHub URL:

https://github.com/tranhoangthienphu/Applied-Data-Science-Capstone/blob/cf92e8c938f8d0ced1b6073477af64fbad6e9471/SpaceX_Machine%20Learning%20Prediction_Part_5.ipynb

# Results

- Exploratory data analysis results: Uncovered relationships between payload mass, booster versions, launch sites, and mission outcomes using bar charts, pie charts, and scatter plots.

- Interactive analytics demo in screenshots: Developed an interactive dashboard (Plotly Dash) and map (Folium) to explore launch data by site, payload range, and success classification.

- Predictive analysis results: Evaluated four models — Logistic Regression, SVM, Decision Tree, and KNN — using GridSearchCV. Decision Tree achieved the best cross-validation accuracy.
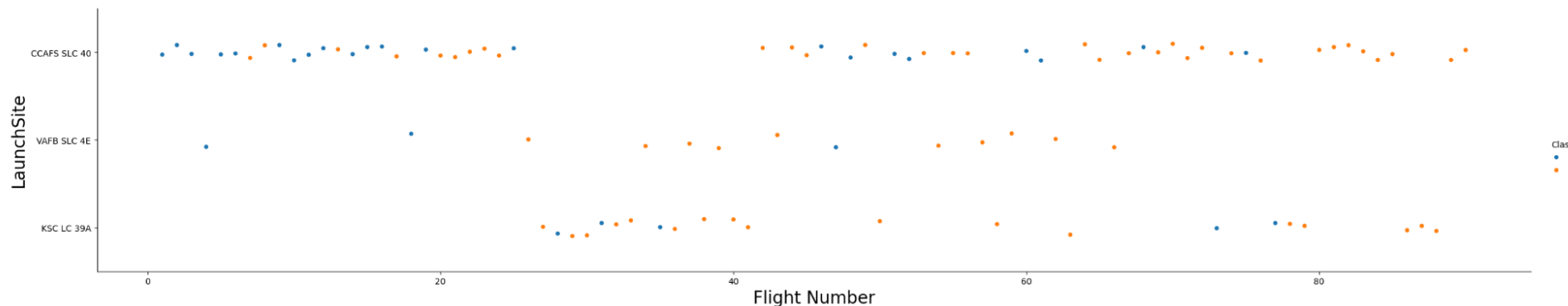
Section 2

# Insights drawn from EDA

# Flight Number vs. Launch Site

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

```
# Plot a scatter point chart with x axis to be Flight Number and y axis to be the launch site, and hue to be the class value
# Plot a scatter point chart with x axis to be Flight Number and y axis to be the launch site, and hue to be the class value
sns.catplot(y="LaunchSite", x="FlightNumber", hue="Class", data=df, aspect = 5)
plt.xlabel("Flight Number",fontsize=20)
plt.ylabel("LaunchSite",fontsize=20)
plt.show()
```
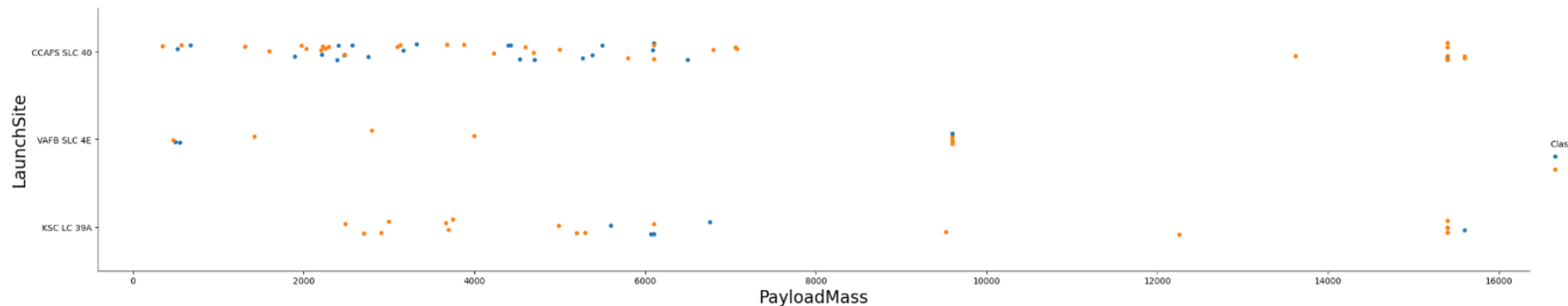


- CCAFS SLC 40 has the highest number of launches, with a balanced mix of successes and failures early on, but mostly successful outcomes in later flights.

- KSC LC 39A shows a high success rate, especially in later flights.

- VAFB SLC 4E has fewer launches, with mixed outcomes.

- Conclusion: Launch success generally improves with higher flight numbers (experience), and success rates vary slightly by launch site.

19

# Payload vs. Launch Site

- Show a scatter plot of Payload vs. Launch Site



```
[6]:  # Plot a scatter point chart with x axis to be Pay Load Mass (kg) and y axis to be the Launch site, and hue to be the class value
      # Plot a scatter point chart with x axis to be Pay Load Mass (kg) and y axis to be the Launch site, and hue to be the class value
      sns.catplot(y="LaunchSite", x="PayloadMass", hue="Class", data=df, aspect = 5)
      plt.xlabel("PayloadMass",fontsize=20)
      plt.ylabel("LaunchSite",fontsize=20)
      plt.show()
```
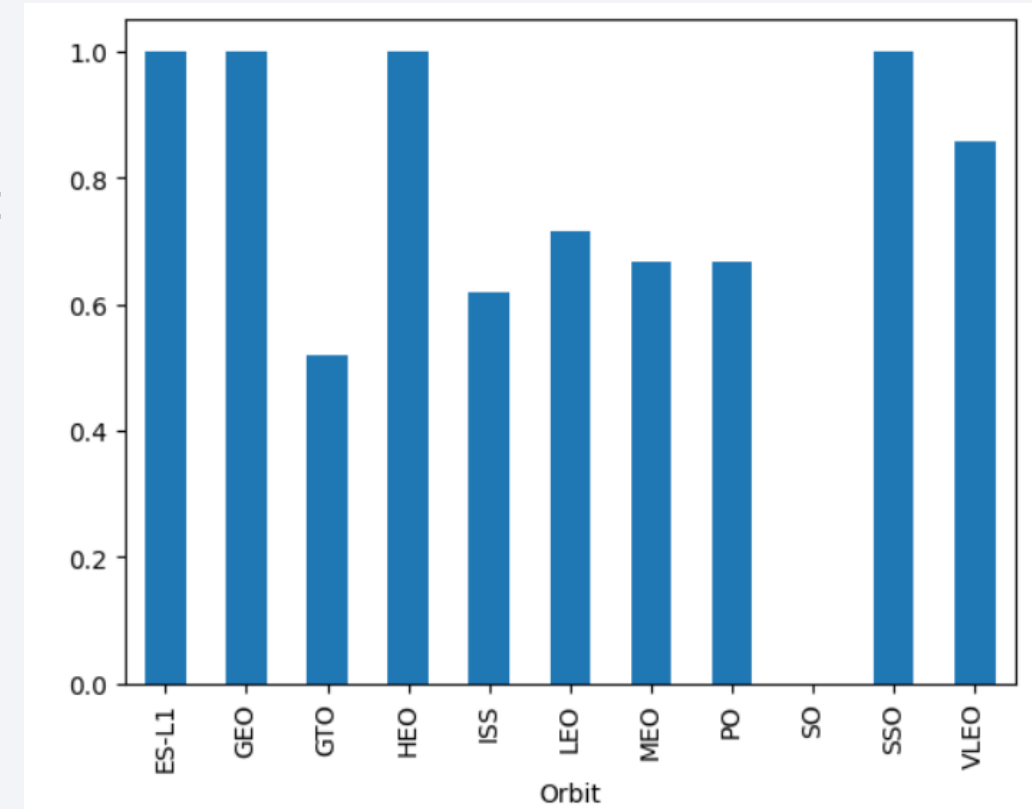
- Most launches with payloads below 6000 kg occurred at CCAFS SLC 40, showing a mix of success and failure.

- High payload launches (over 10,000 kg) are more common at KSC LC 39A, and most of them were successful.

- VAFB SLC 4E had fewer launches and lower payload masses, with varying outcomes.

- **Conclusion:** Higher payloads, especially from **KSC LC 39A**, tend to have higher success rates, suggesting advanced readiness at that site.

# Success Rate vs. Orbit Type

- Orbits such as ES-L1, GEO, HEO, and SSO have a 100% success rate, indicating highly reliable missions to these destinations.

- GTO (Geostationary Transfer Orbit) has the lowest success rate, around 50%, suggesting more technical challenges or complexity.

- Common orbits like LEO (Low Earth Orbit) and MEO (Medium Earth Orbit) have moderate success rates between 65%–75%.
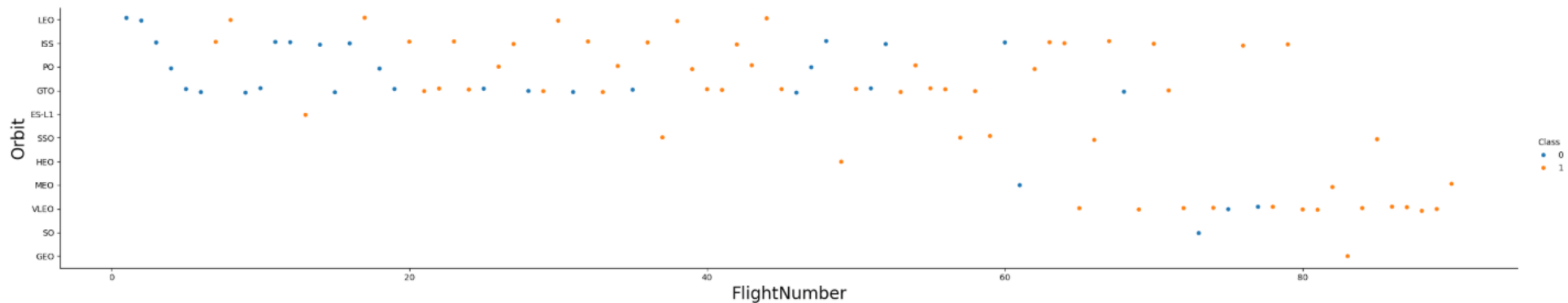
Conclusion: The chart highlights that while most orbit types maintain high success rates, special attention may be needed for missions targeting GTO, where failure rates are noticeably higher.

# Flight Number vs. Orbit Type

- The success of missions in LEO seems to benefit from accumulated launch experience, showing a trend of increasing reliability. In contrast, GTO missions show no strong correlation with flight history, suggesting intrinsic challenges associated with that orbit



```
[8]:  # Plot a scatter point chart with x axis to be FlightNumber and y axis to be the Orbit, and hue to be the class value
      sns.catplot(y="Orbit", x="FlightNumber", hue="Class", data=df, aspect = 5)
      plt.xlabel("FlightNumber",fontsize=20)
      plt.ylabel("Orbit",fontsize=20)
      plt.show()
```
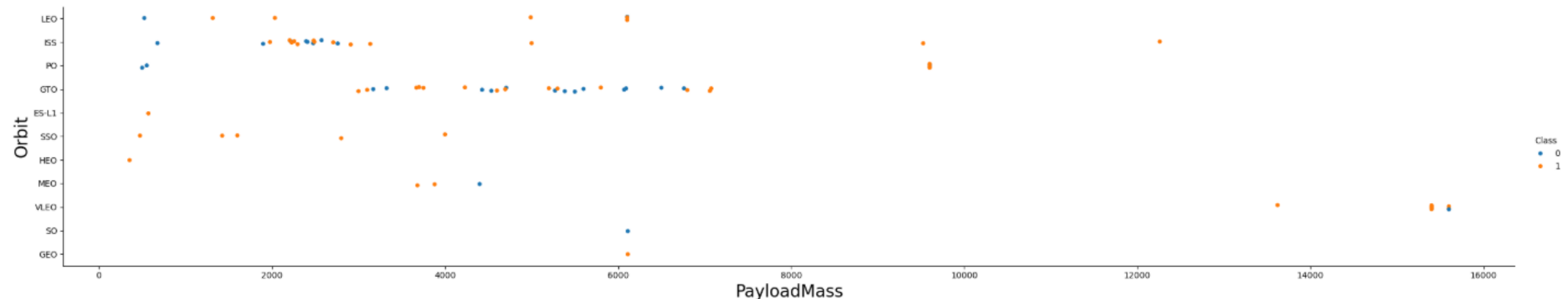
You can observe that in the LEO orbit, success seems to be related to the number of flights. Conversely, in the GTO orbit, there appears to be no relationship between flight number and success.

# Payload vs. Orbit Type

- While SpaceX missions with heavy payloads tend to show higher success rates for LEO, ISS, and Polar orbits, missions targeting GTO display a mixed success pattern, reflecting greater operational complexity or risks.



```
[9]:  # Plot a scatter point chart with x axis to be Payload and y axis to be the Orbit, and hue to be the class value
      sns.catplot(y="Orbit", x="PayloadMass", hue="Class", data=df, aspect = 5)
      plt.xlabel("PayloadMass",fontsize=20)
      plt.ylabel("Orbit",fontsize=20)
      plt.show()
```

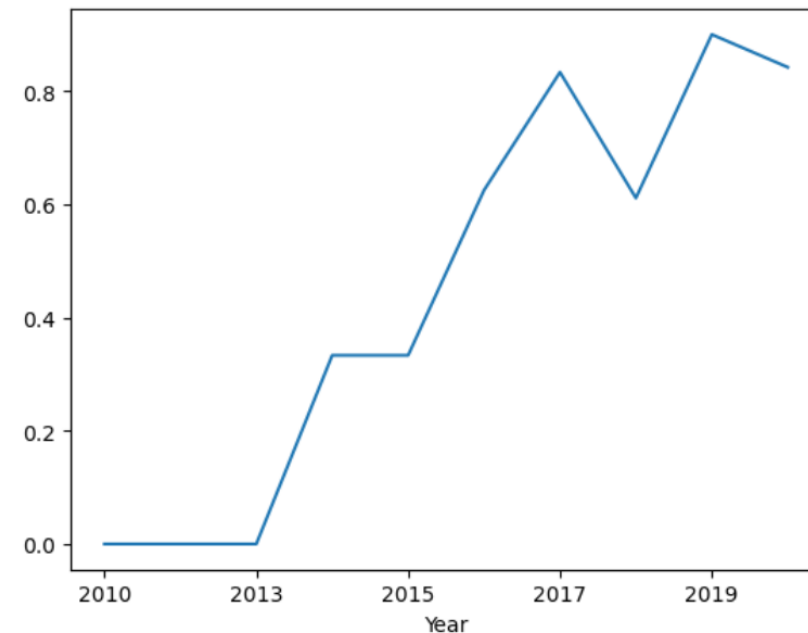With heavy payloads the successful landing or positive landing rate are more for Polar,LEO and ISS.

However, for GTO, it's difficult to distinguish between successful and unsuccessful landings as both outcomes are present.

23

# Launch Success Yearly Trend

- Since 2013, SpaceX has shown a strong upward trend in mission success rate, peaking around 2019. This reflects consistent improvements in reliability and operational efficiency over time

```
[11]:  # Plot a line chart with x axis to be the extracted year and y axis to be the success rate
       temp_df = df.copy()
       temp_df['Year'] = year
       temp_df.groupby('Year')['Class'].mean().plot()

[11]:  <AxesSubplot:xlabel='Year'>
```



you can observe that the sucess rate since 2013 kept increasing till 2020

# All Launch Site Names

Distinct Launch Sites:

- CCAFS LC-40 – Cape Canaveral Air Force Station, Launch Complex 40

- VAFB SLC-4E – Vandenberg Air Force Base, Space Launch Complex 4E

- KSC LC-39A – Kennedy Space Center, Launch Complex 39A

- CCAFS SLC-40 – (Appears to be a duplicate or naming inconsistency of CCAFS LC-40)

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

```
In [24]:  %%sql
          SELECT DISTINCT "Launch_Site"
          FROM SPACEXTBL;
```

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

Out[24]:

| Launch_Site |
| --- |
| CCAFS LC-40 |
| VAFB SLC-4E |
| KSC LC-39A |
| CCAFS SLC-40 |

# Launch Site Names Begin with 'CCA'

- NASA's Commercial Resupply Services (CRS) missions launched by SpaceX delivered a total of 45,596 kg of payload to their destinations. This highlights the significant contribution of SpaceX to NASA's supply chain to the International Space Station (ISS).

## Task 2

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

```sql
[11]: %%sql
SELECT *
FROM SPACEXTBL
WHERE "Launch_Site" LIKE 'CCA%'
LIMIT 5;
```

 * sqlite:///my_data1.db
Done.

[11]:

| 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

- Calculate the total payload carried by boosters from NASA

- Present your query result with a short explanation here

## Task 3

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

```sql
[12]: %%sql
SELECT SUM("PAYLOAD_MASS__KG_") AS "Total_Payload_Mass_KG"
FROM SPACEXTBL
WHERE "Customer" = 'NASA (CRS)';
```

 * sqlite:///my_data1.db
Done.

[12]: **Total_Payload_Mass_KG**

45596

# Average Payload Mass by F9 v1.1

- The Falcon 9 version 1.1 booster typically carried a payload mass of around 2928.4 kg per launch, reflecting its operational capacity before newer versions like F9 FT and F9 B5 were introduced with higher performance.

## Task 4

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

```sql
[13]:  %%sql
       SELECT AVG("PAYLOAD_MASS__KG_") AS "Average_Payload_Mass_KG"
       FROM SPACEXTBL
       WHERE "Booster_Version" = 'F9 v1.1';
```

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

[13]: **Average_Payload_Mass_KG**

2928.4

# First Successful Ground Landing Date

- SpaceX achieved its first ground pad landing success on December 22, 2015. This marked a major milestone in reusable rocket technology, showcasing the company's ability to recover and potentially reuse rocket boosters landed on solid ground.

# Successful Drone Ship Landing with Payload between 4000 and 6000

- These four Falcon 9 Full Thrust (FT) boosters successfully landed on drone ships while carrying moderate to heavy payloads (between 4000–6000 kg). This highlights the reliability and reusability of these booster versions for mid-weight orbital missions.

# Total Number of Successful and Failure Mission Outcomes

- SpaceX missions have a very high success rate, with over 97% of missions marked successful. This highlights the reliability and maturity of their launch systems. However, the slight duplication of "Success" suggests a potential inconsistency in labeling, which could benefit from data cleaning.

## Task 7

List the total number of successful and failure mission outcomes

```sql
[18]: %%sql
SELECT Mission_Outcome, COUNT(*) AS Total
FROM SPACEXTBL
GROUP BY Mission_Outcome;
```

 * sqlite:///my_data1.db
Done.

[18]:

| Mission_Outcome | Total |
|---|---|
| Failure (in flight) | 1 |
| Success | 98 |
| Success | 1 |
| Success (payload status unclear) | 1 |

# Boosters Carried Maximum Payload

- All listed boosters are of type **F9 B5**, which is the **Block 5 variant** of Falcon 9.

- Examples: F9 B5 B1048.4, B1049.4, B1051.3, etc.

- This indicates that the **Block 5** version is responsible for **heaviest payload deliveries**, showcasing its enhanced capacity.

Conclusion: This insight supports the conclusion that newer booster versions (like Block 5) play a crucial role in high-performance missions.

List all the booster_versions that have carried the maximum payload mass, using a subquery with a suitable aggregate function.

```sql
[21]: %%sql
SELECT Booster_Version
FROM SPACEXTBL
WHERE Payload_Mass__kg_ = (
    SELECT MAX(Payload_Mass__kg_)
    FROM SPACEXTBL
);
```

 * sqlite:///my_data1.db
Done.

[21]: | 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 |

# 2015 Launch Records

- **January 2015**: Failure using **F9 v1.1 B1012** at **CCAFS LC-40**

- **April 2015**: Failure using **F9 v1.1 B1015** at **CCAFS LC-40**

## Task 9

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

 * sqlite:///my_data1.db
Done.

[22]:

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

33

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

- **Most frequent outcome**: **No attempt** (10 launches) – indicating many missions did not include a landing attempt.Present your query result with a short explanation here

- **Success (drone ship)** and **Failure (drone ship)** both occurred **5 times** each, showing a balanced but challenging performance on drone ships.

- **Success (ground pad)** was recorded **3 times**, same as **Controlled (ocean)** – both are less frequent but significant.

- Rare outcomes include **Failure (parachute),** **Uncontrolled (ocean)**, and **Precluded (drone ship)**.

Conclusion: This result reflects early experimentation and progress in SpaceX's landing strategies during the specified period.

### Task 10

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.

```
[23]: %%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;
```

 * sqlite:///my_data1.db
Done.

[23]:

| 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 |

34

Section 3

# Launch Sites
# Proximities Analysis

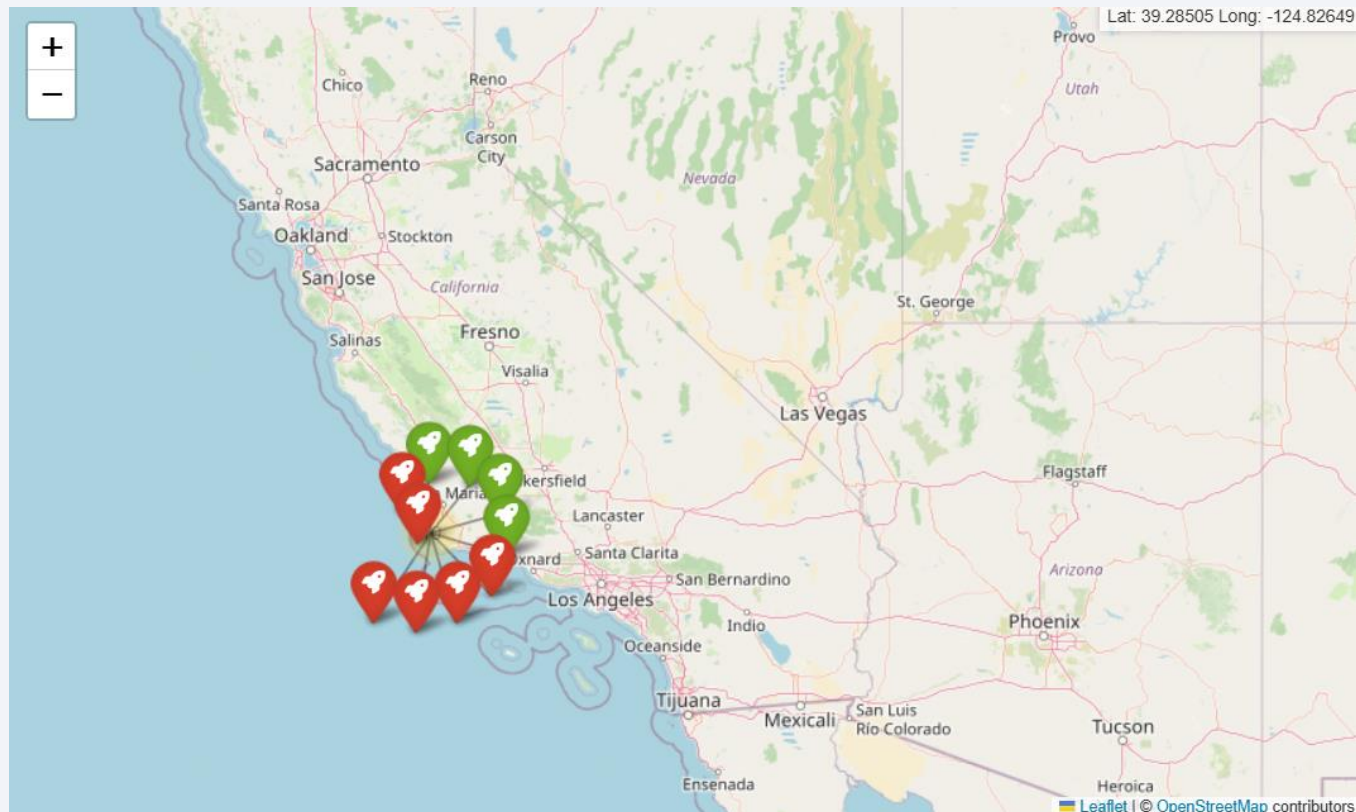# <Folium Map Screenshot 1>

- It can be seen that the map shows the concentration locations of the launches, one launch point has 10 launches, the other has 46 launches.

# <Folium Map Screenshot 2>
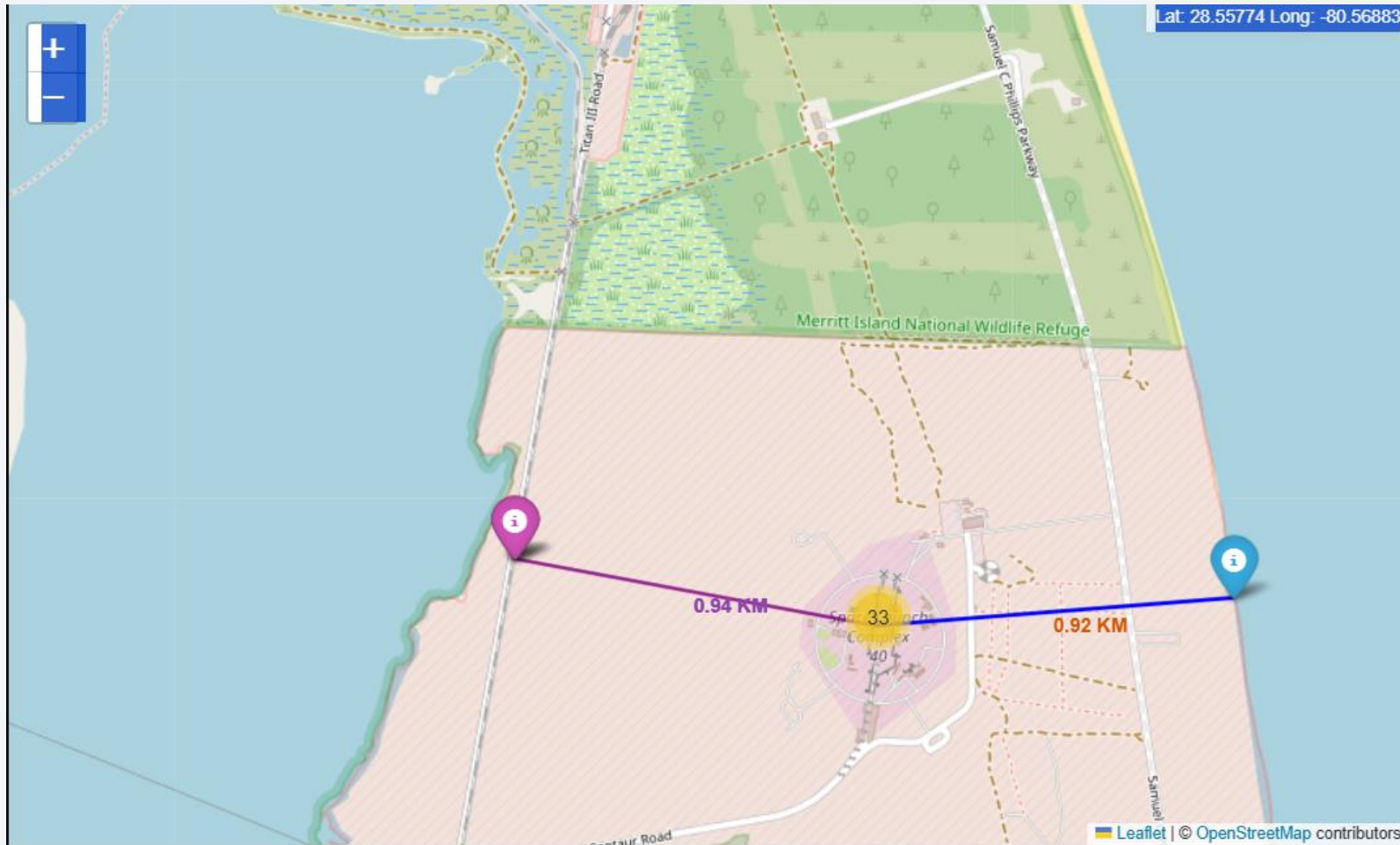
- launsite location includes launched points, green points are success, red points are failure

# <Folium Map Screenshot 3>

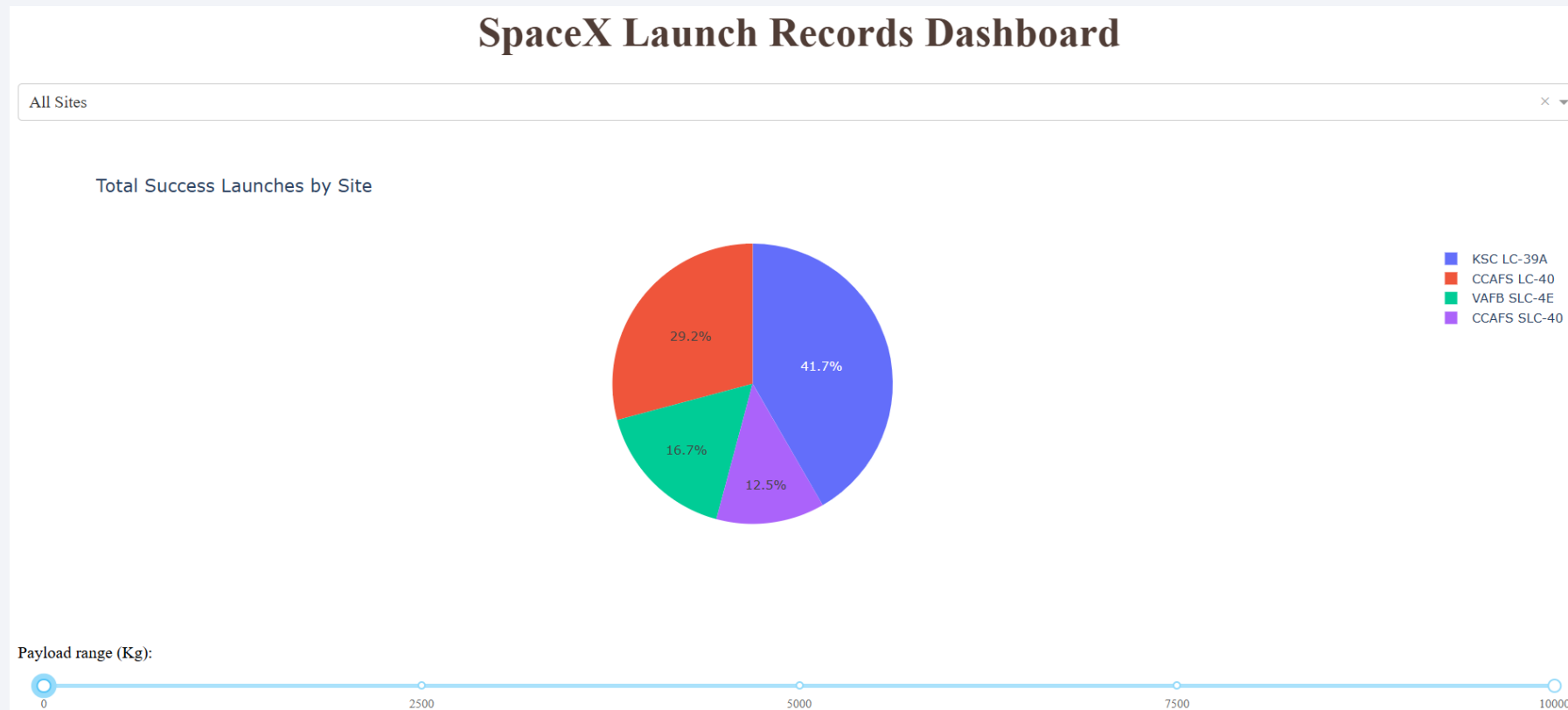- The launch site is 0.94 km from the nearest highway (purple line)

Section 4

# Build a Dashboard
# with Plotly Dash

# \<Dashboard Screenshot 1\>

- KSC LC-39A is the top-performing launch site in terms of successful missions (41.7%)
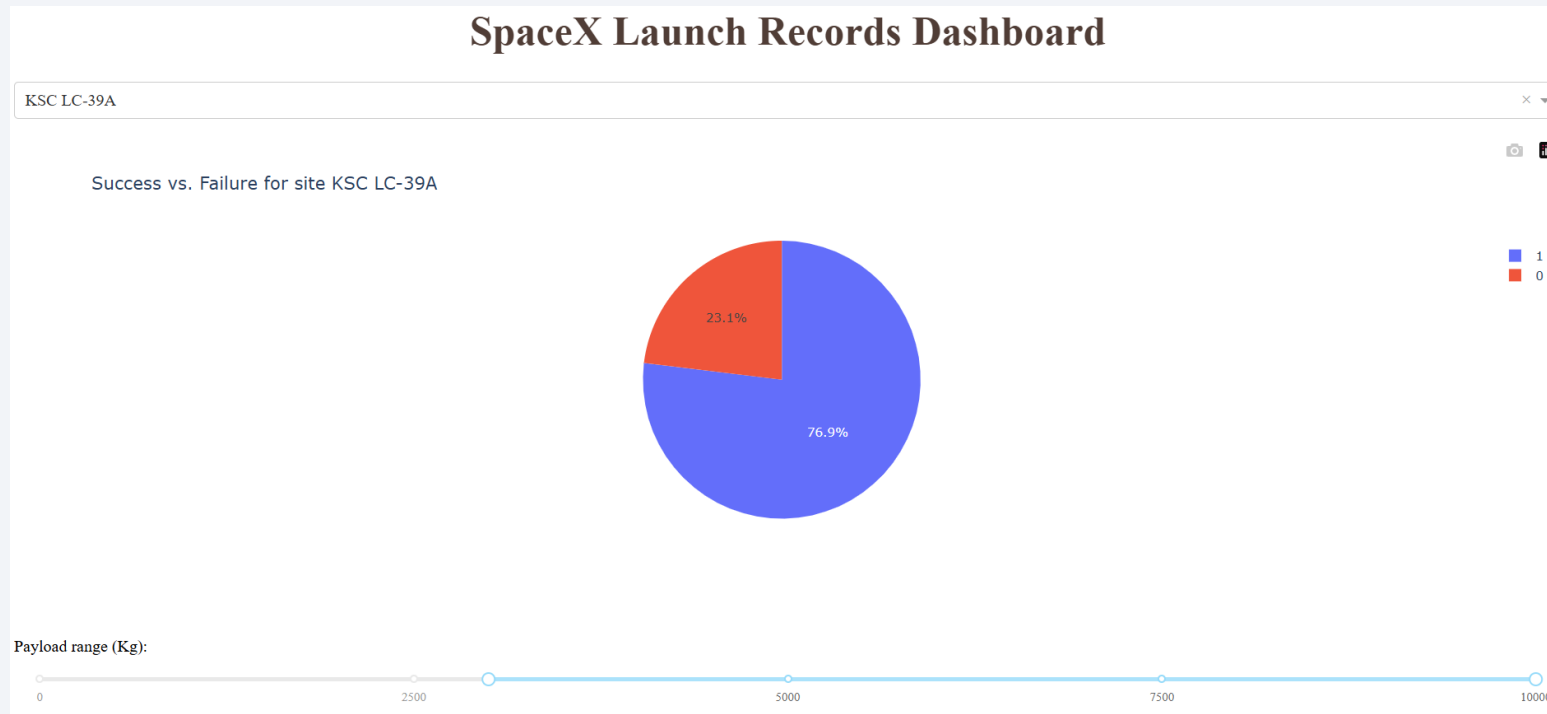
# \<Dashboard Screenshot 2\>

This pie chart from the **SpaceX Launch Records Dashboard** displays the **success vs. failure rate** for launches at **KSC LC-39A**:

- 76.9% of launches were successful (class = 1, shown in blue).

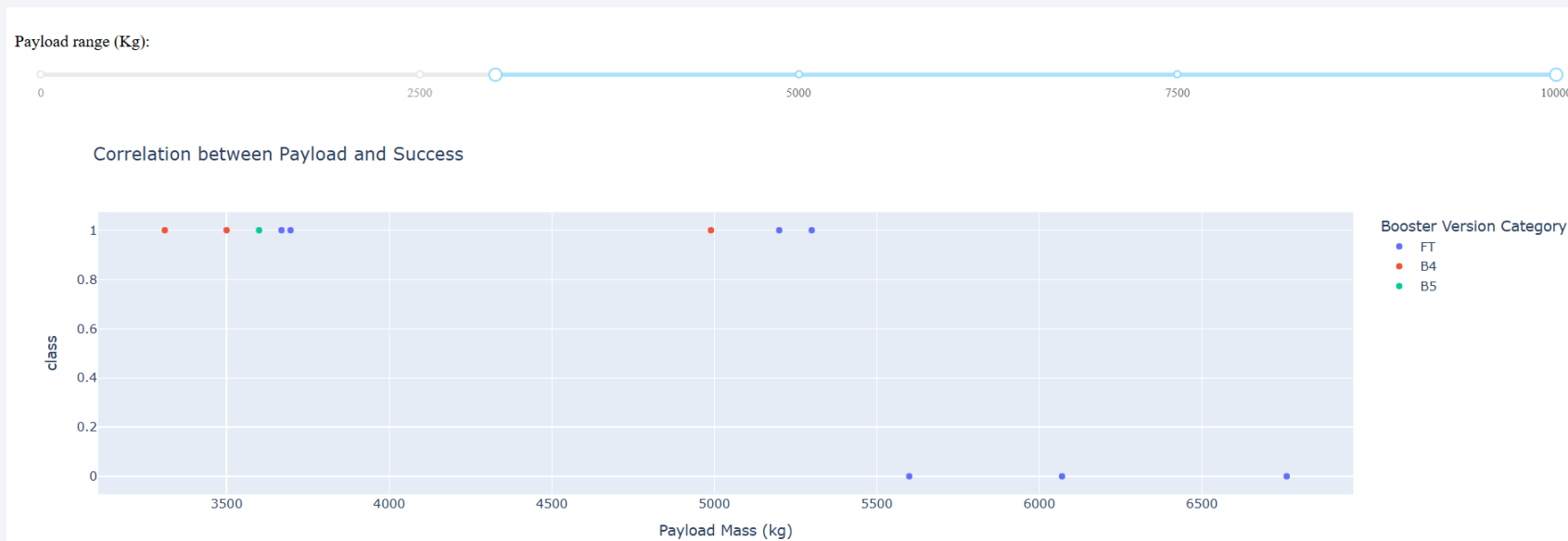- 23.1% of launches failed (class = 0, shown in red).

KSC LC-39A demonstrates a high launch success rate, making it one of the most reliable launch sites.



**SpaceX Launch Records Dashboard**

KSC LC-39A

Success vs. Failure for site KSC LC-39A

23.1%

76.9%

1
0

Payload range (Kg):

0          2500          5000          7500          10000

41

# &lt;Dashboard Screenshot 3&gt;

- Most missions with payloads between 3500–5000 kg were successful (class = 1).

- Failures (class = 0) tend to appear more as the payload mass increases beyond 5500 kg.

- Booster version FT appears more frequently and shows both successes and failures.

Conclusion: Lighter payloads (especially below ~5000 kg) are **more likely to succeed**, regardless of booster type.
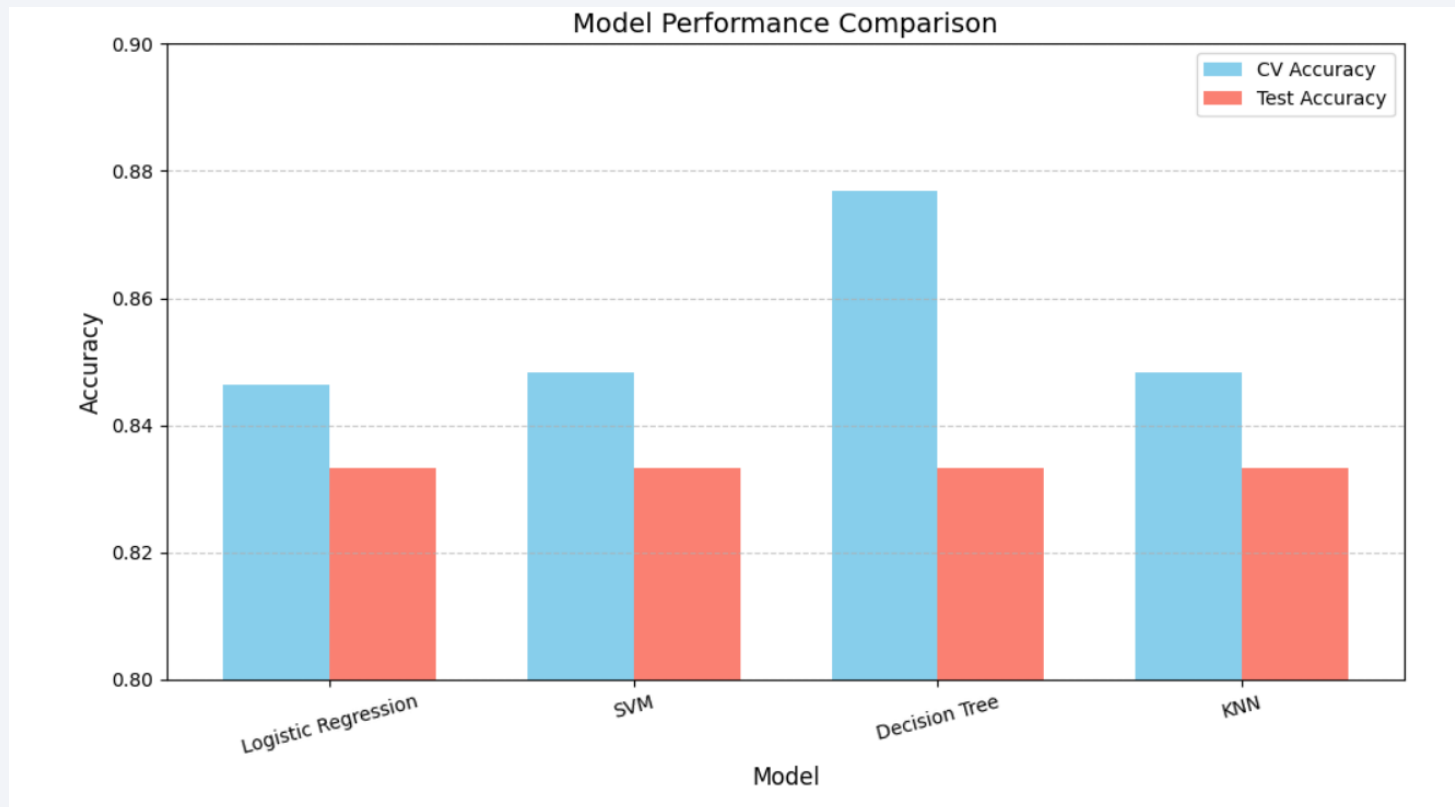
Section 5

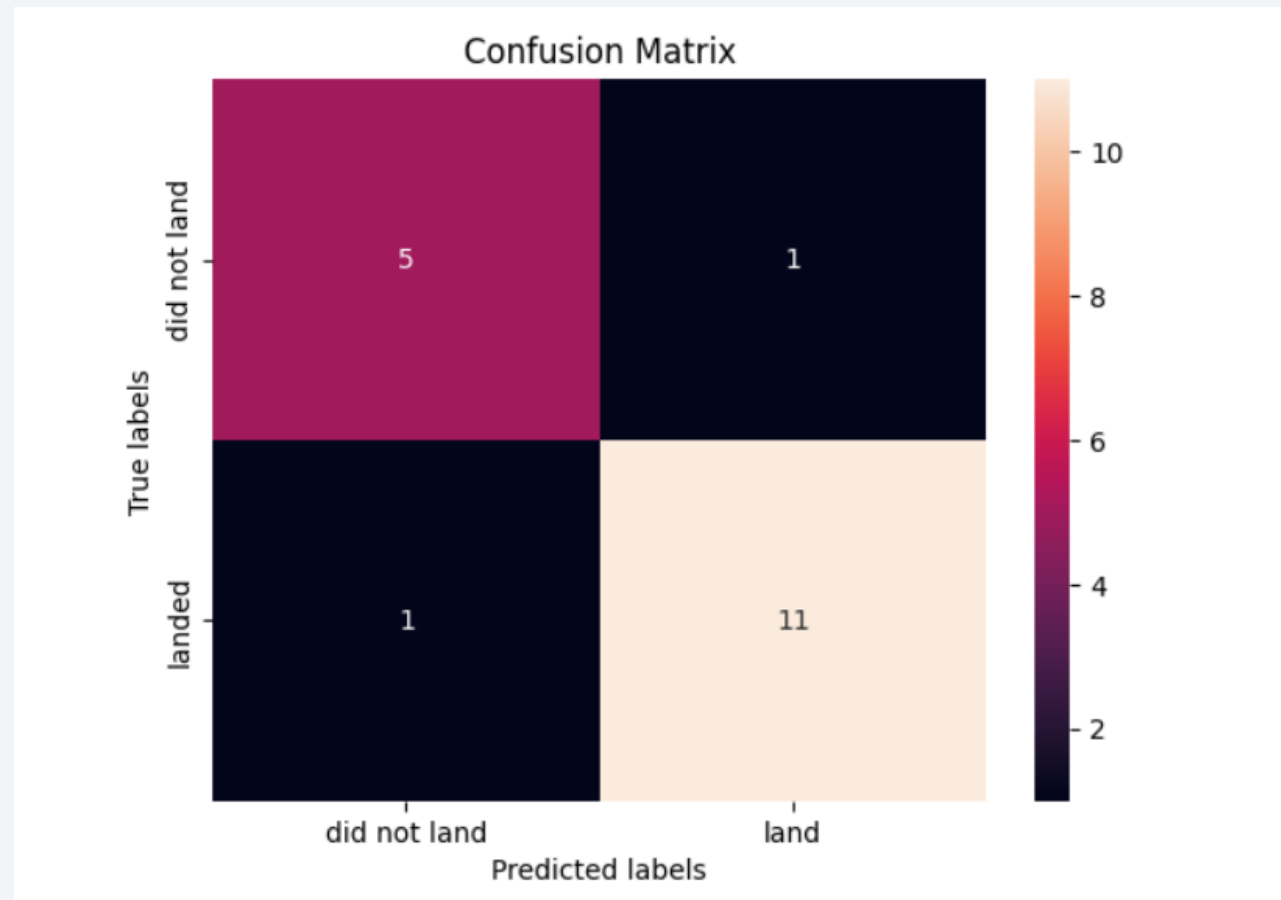# Predictive Analysis (Classification)

# Classification Accuracy

- Although Decision Tree had the best training performance, **all models generalized similarly** on test data. Hence, **Decision Tree** was selected for its balance of performance and interpretability.

# Confusion Matrix

- The model performed well, correctly predicting 16 out of 18 cases. The balanced low number of false predictions indicates high accuracy and good generalization.



Confusion Matrix

# Conclusions

- Multiple data sources were explored to iteratively refine insights and support data-driven conclusions.

- Among all launch sites, KSC LC-39A emerged as the most reliable and successful.

- Launch missions carrying payloads exceeding 7,000 kg demonstrated a lower risk profile.

- While the majority of missions achieved success, the success rate of landings has shown consistent improvement over time—reflecting advancements in rocket design and operational procedures.

- The Decision Tree Classifier proved to be the most effective model for predicting landing outcomes, offering potential for enhanced mission planning and profitability.

# 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!