



IBM Developer
SKILLS NETWORK

Winning Space Race with Data Science

Muhammad Sohail
February 15, 2025



Outline

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

Executive Summary

Summary of Methodologies

1. Data Collection: Accessed SpaceX launch data via API and web scraped records from Wikipedia.

2. Data Cleaning & Preparation:

- Cleaned and formatted the data.
- Stored data in Db2 database and performed SQL queries.
- Conducted exploratory data analysis.

3. Feature Engineering: Created new features and standardized the data.

4. Interactive Visualizations:

- Mapped launch sites and success rates using Folium.
- Built an interactive dashboard with Plotly Dash.

5. Model Building & Evaluation:

- Implemented SVM, Decision Trees, and K-Nearest Neighbors.
- Tuned hyperparameters with GridSearchCV.
- Evaluated models using test data accuracy.

Summary of Results

1. Data Insights:

- Identified factors influencing Falcon 9 first stage landings.
- Visualized geographical patterns and success rates.

2. Model Performance:

- SVM : 83.33% accuracy.
- K-Nearest Neighbors : 83.33 % accuracy.
- Decision Tree: 83.33 % accuracy.

3. Key Findings:

- Launch site and payload mass impact landing success.
- Decision Tree model is the most effective predictor.

Introduction

Project Background and Context:

In this capstone project, we aim to predict the successful landing of the Falcon 9 first stage. SpaceX advertises rocket launches at a significantly lower cost compared to other providers, largely due to their ability to reuse the first stage of the rocket. By accurately predicting landing success, we can estimate launch costs and provide valuable insights for companies bidding against SpaceX.

Problems We Want to Find Answers To:

- What factors influence the successful landing of the Falcon 9 first stage?
- How can we accurately predict the landing outcome using machine learning models?
- Which machine learning model performs best in predicting the landing success?

Section 1

Methodology

Methodology

Executive Summary: This project employs a comprehensive approach to predict the successful landing of the Falcon 9 first stage, incorporating data collection, processing, exploratory analysis, interactive visualizations, and predictive modeling.

Data Collection Methodology: Data was sourced from the SpaceX API, which provided detailed records of Falcon 9 launches, including launch dates, sites, payloads, and outcomes.

Perform Data Wrangling: Data cleaning involved handling missing values, standardizing formats, and ensuring consistency. Key features were extracted and new features engineered to enrich the dataset.

Perform Exploratory Data Analysis (EDA) Using Visualization and SQL:

- Visualized launch success rates, payloads, and launch sites using Matplotlib and Seaborn.
- Executed SQL queries to derive insights and answer specific questions regarding the dataset.

Methodology

Perform Interactive Visual Analytics Using Folium and Plotly Dash:

- Used Folium to create interactive maps displaying launch sites and outcomes.
- Developed a Plotly Dash application with interactive components like dropdowns and sliders to analyze launch success rates and payload ranges.

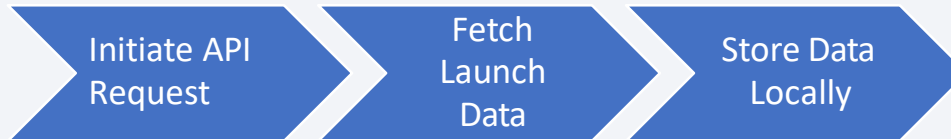
Perform Predictive Analysis Using Classification Models:

- Built and evaluated various classification models including Logistic Regression, SVM, KNN, and Decision Trees.
- Employed GridSearchCV for hyperparameter tuning.
- Evaluated models based on accuracy, and identified the best performing model for predicting landing success.

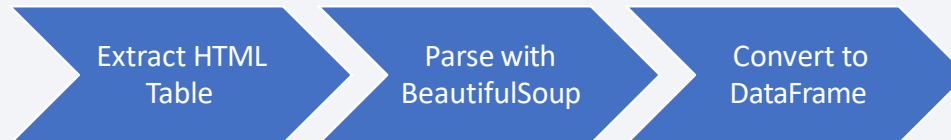
Github URL: <https://github.com/theitlink/assignment.git>

Data Collection

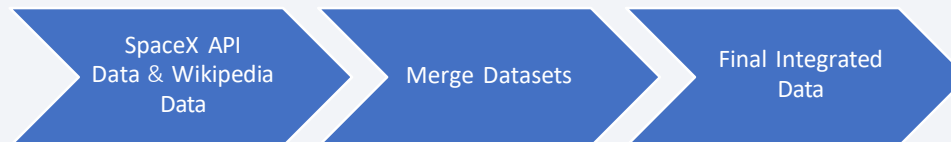
- Step 1: SpaceX API Request



- Step 2: Web Scraping Wikipedia



- Step 3: Data Integration



Data sets were collected from:

Space X API (<https://api.spacexdata.com/v4/rockets/>)

Wikipedia (https://en.wikipedia.org/wiki/List_of_Falcon_9_and_Falcon_Heavy_launches), using web scraping technics.

Data Collection – SpaceX API

Step 1: Initiate API Request

- Use Python's `requests` library to connect to the SpaceX API.
- Endpoint: `https://api.spacexdata.com/v4/launches`

Step 2: Parse API Response

- Convert API response from JSON to a Python dictionary.
- Extract relevant fields: launch date, launch site, payload mass, rocket type, outcome.

Step 3: Store Data Locally

- Save extracted data into a pandas DataFrame.
- Store the DataFrame locally for further processing.

GitHub URL: 1. [jupyter-labs-spacex-data-collection-api.ipynb](#)



Data Collection - Scraping

Step 1: Initiate Web Scraping

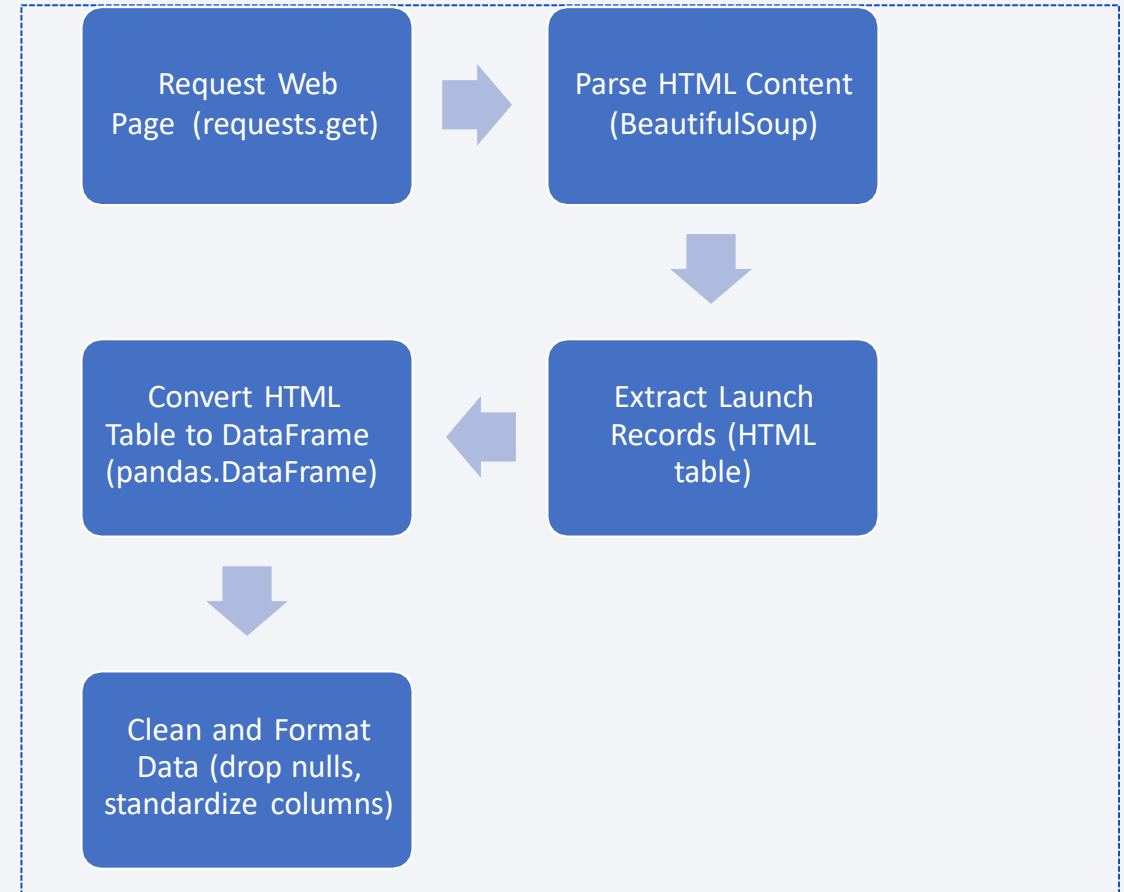
- Use Python's `requests` library to fetch the HTML content of the Wikipedia page.
- Target URL:
`https://en.wikipedia.org/wiki/List_of_Falcon_9_and_Falcon_Heavy_launches`

Step 2: Parse HTML Content

- Use `BeautifulSoup` to parse the HTML content.
- Extract the HTML table containing Falcon 9 launch records.

Step 3: Convert to DataFrame

- Convert the extracted HTML table into a pandas DataFrame.
- Clean and format the DataFrame, ensuring data consistency.



Data Wrangling

Overview: Data wrangling involves cleaning, transforming, and organizing raw data into a structured format suitable for analysis.

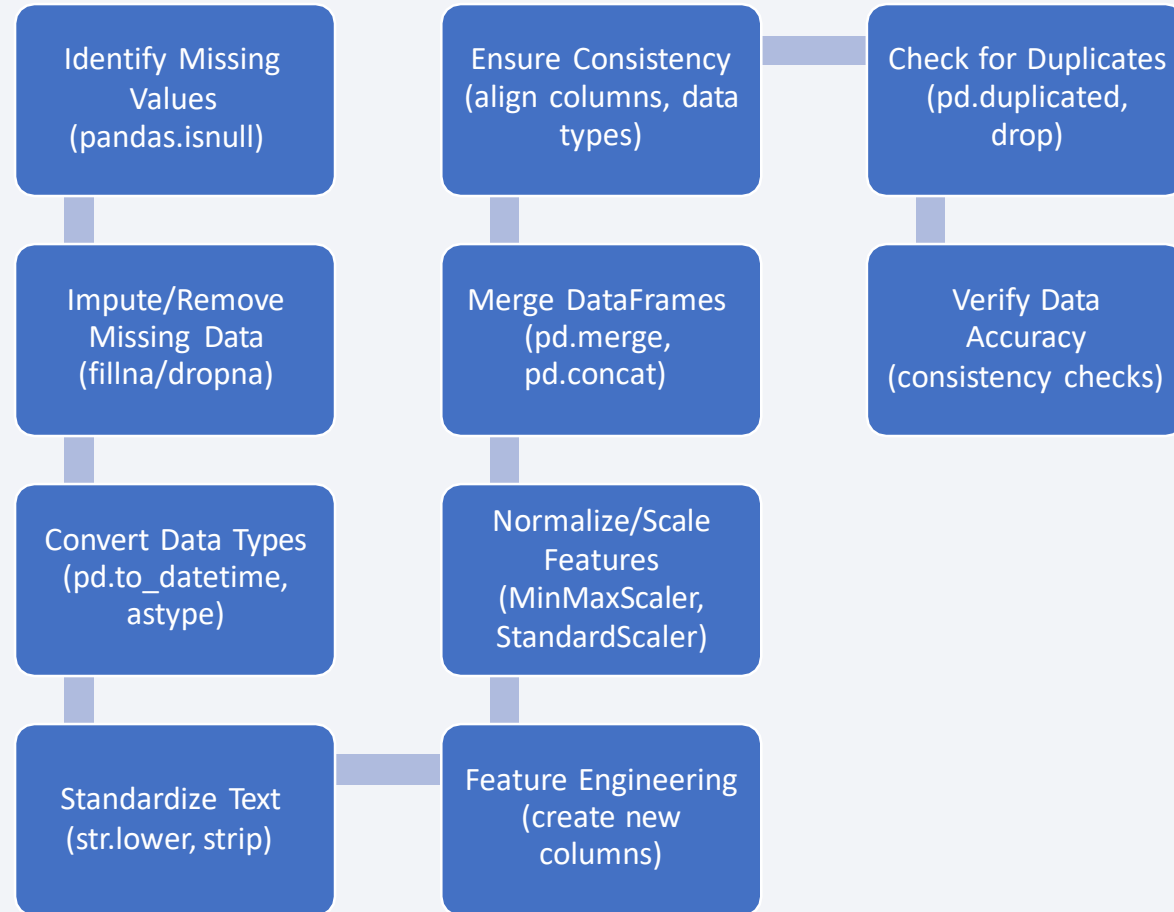
- Step 1: Data Cleaning
 - Identify and fill or remove missing values in the dataset.
 - Use appropriate imputation techniques or drop rows/columns with excessive missing data.
- Step 2: Data Transformation
 - Convert data types to appropriate formats (e.g., date-time, numerical).
 - Standardize text (e.g., lowercase, remove whitespace).
 - Create new features from existing data (e.g., extract year from date).
 - Normalize/scale numerical features to ensure consistency.

Data Wrangling

- Step 3: Data Integration
 - Merge datasets collected from different sources (API, web scraping) into a single cohesive dataset.
 - Ensure consistent column names and data formats across datasets.
- Step 4: Data Validation
 - Check for duplicate records and remove them.
 - Verify the accuracy and consistency of data entries.

GitHub URL: [2. labs-jupyter-spacex-Data wrangling.ipynb](#)

Data Wrangling



EDA with Data Visualization

Overview:

Exploratory Data Analysis (EDA) involves visually exploring and summarizing the main characteristics of a dataset. The goal is to understand the data's distribution, identify patterns, and uncover relationships between variables.

Charts Plotted:

1. Histograms:

- **Purpose:** Used to visualize the distribution of numerical variables such as launch success rates, payload mass, and flight number.
- **Why:** Helps in understanding the spread and central tendency of the data, identifying outliers, and assessing data skewness.

2. Bar Charts:

- **Purpose:** Used to compare categorical variables such as launch outcomes (success/failure) across different categories like launch sites or rocket types.
- **Why:** Provides a clear comparison of frequencies or proportions within categorical data, highlighting patterns or trends.

3. Line Charts:

- **Purpose:** Used to track trends over time, such as the success rate of Falcon 9 launches across different years.
- **Why:** Reveals temporal patterns and helps in understanding performance trends or changes over specific periods.

EDA with Data Visualization

4. Scatter Plots:

- **Purpose:** Used to explore relationships between two numerical variables, such as payload mass vs. launch success.
- **Why:** Identifies correlations or dependencies between variables, visualizing how one variable changes concerning another.

5. Heatmaps:

- **Purpose:** Used to visualize correlation matrices between multiple numerical variables.
- **Why:** Helps in identifying strong correlations (positive or negative) between variables, aiding feature selection or understanding multicollinearity.

6. Box Plots:

- **Purpose:** Used to display the distribution of numerical data through their quartiles.
- **Why:** Visualizes the spread and skewness of data, highlighting outliers and comparing distributions across different categories.

Github URL: [3. edadataviz.ipynb](#)

EDA with SQL

Aggregate Queries:

- Calculated total number of launches.
- Counted successful and failed launches.
- Calculated success rates by launch site and rocket type.

Join Queries:

- Joined tables to link launch records with additional data (e.g., rocket details).
- Combined datasets for comprehensive analysis.

Filtering Queries:

- Filtered data to focus on specific launch outcomes (success/failure).
- Applied conditions to extract launches based on criteria like launch date or rocket configuration.

Sorting Queries:

- Sorted data to identify trends or outliers.
- Ordered launches by date or success rate for analysis.

Subqueries:

- Nested queries to calculate derived metrics (e.g., average payload mass per launch site).
- Subqueries used to perform detailed analysis within larger datasets.

GitHub URL: [4. jupyter-labs-eda-sql-coursera_sqlite.ipynb](#)

Build an Interactive Map with Folium

Map Objects Created

Markers:

- Placed markers to indicate launch sites on the map.
- Each marker represents a specific geographical location where SpaceX launches have occurred.

Circles:

- Added circles around launch sites to visually represent proximity zones.
- Circles help visualize the areas around launch sites that might influence operational decisions.

Lines:

- Drew lines to connect launch sites with their proximities or other relevant locations.
- Lines provide spatial context and connections between different points of interest related to launches.

Reasons for Adding Objects

Markers:

- To pinpoint exact launch locations for spatial reference.
- Helps users identify where SpaceX has conducted launches geographically.

Circles:

- Illustrates the potential impact zones around launch sites.
- Provides a visual representation of safety perimeters or operational boundaries.

Lines:

- Shows connections or relationships between launch sites and relevant features.
- Enhances understanding of spatial relationships and dependencies.

Github URL: [5. lab_jupyter_launch_site_location.ipynb](#)

Build a Dashboard with Plotly Dash

Plots/Graphs Added

Success Pie Chart:

- Displays the distribution of successful and failed launches.
- Helps visualize the overall success rate and performance trends.

Success-Payload Scatter Plot:

- Shows the relationship between payload mass and launch success.
- Allows users to explore how payload mass influences mission outcomes.

Github URL: `6. spacex_dash_app.py`

Interactions Added

Launch Site Dropdown:

- Enables users to select specific launch sites for analysis.
- Facilitates filtering and focused exploration based on geographical locations.

Range Slider for Payload:

- Allows users to adjust payload mass ranges dynamically.
- Offers flexibility in examining launch success concerning payload mass variations.

Build a Dashboard with Plotly Dash

Reasons:

Success Pie Chart:

- Provides a quick overview of mission success rates.
- Essential for stakeholders to understand overall performance metrics at a glance.

Success-Payload Scatter Plot:

- Helps identify correlations between payload characteristics and launch outcomes.
- Supports decision-making processes related to payload planning and operational strategies.

Launch Site Dropdown:

- Enhances user experience by focusing analysis on specific launch locations.
- Allows for regional insights and comparisons across different launch sites.

Range Slider for Payload:

- Offers interactive exploration of how payload mass affects mission success.
- Enables detailed analysis and insights into payload-related performance factors.

Predictive Analysis (Classification)

1. Data Preprocessing:

- Standardized features to ensure all variables contribute equally.
- Split data into training and test sets for model validation.

2. Model Selection:

- Explored multiple classification algorithms: SVM, Decision Trees, and K-Nearest Neighbors (KNN).
- Chose algorithms suitable for binary classification tasks based on project requirements.

3. Hyperparameter Tuning:

- Used GridSearchCV to systematically search for optimal hyperparameters.
- Tuned parameters such as C (SVM), max_depth (Decision Trees), and n_neighbors (KNN).

Github URL:

7. SpaceX_Machine_Learning_Prediction_Part_5.ipynb

4. Model Evaluation:

- Evaluated models using cross-validation techniques to ensure robustness and generalizability.
- Utilized metrics like accuracy, precision, recall, and F1-score to assess model performance.

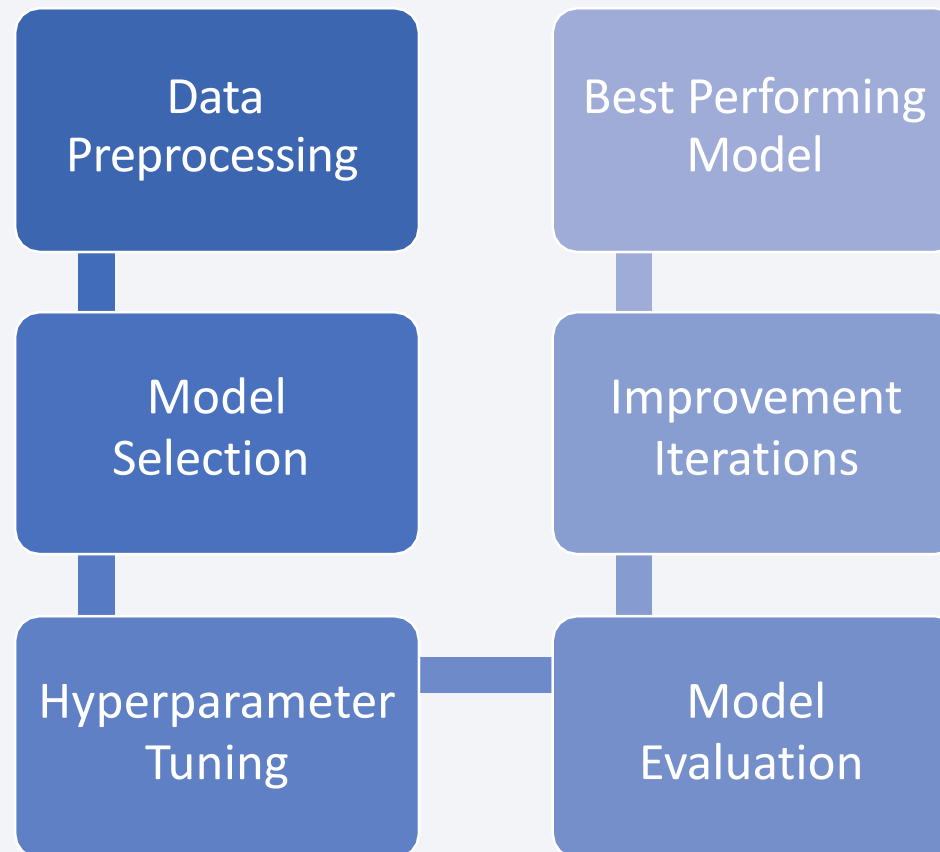
5. Improvement Iterations:

- Iteratively adjusted models based on insights from validation results.
- Fine-tuned hyperparameters to maximize predictive accuracy and reliability.

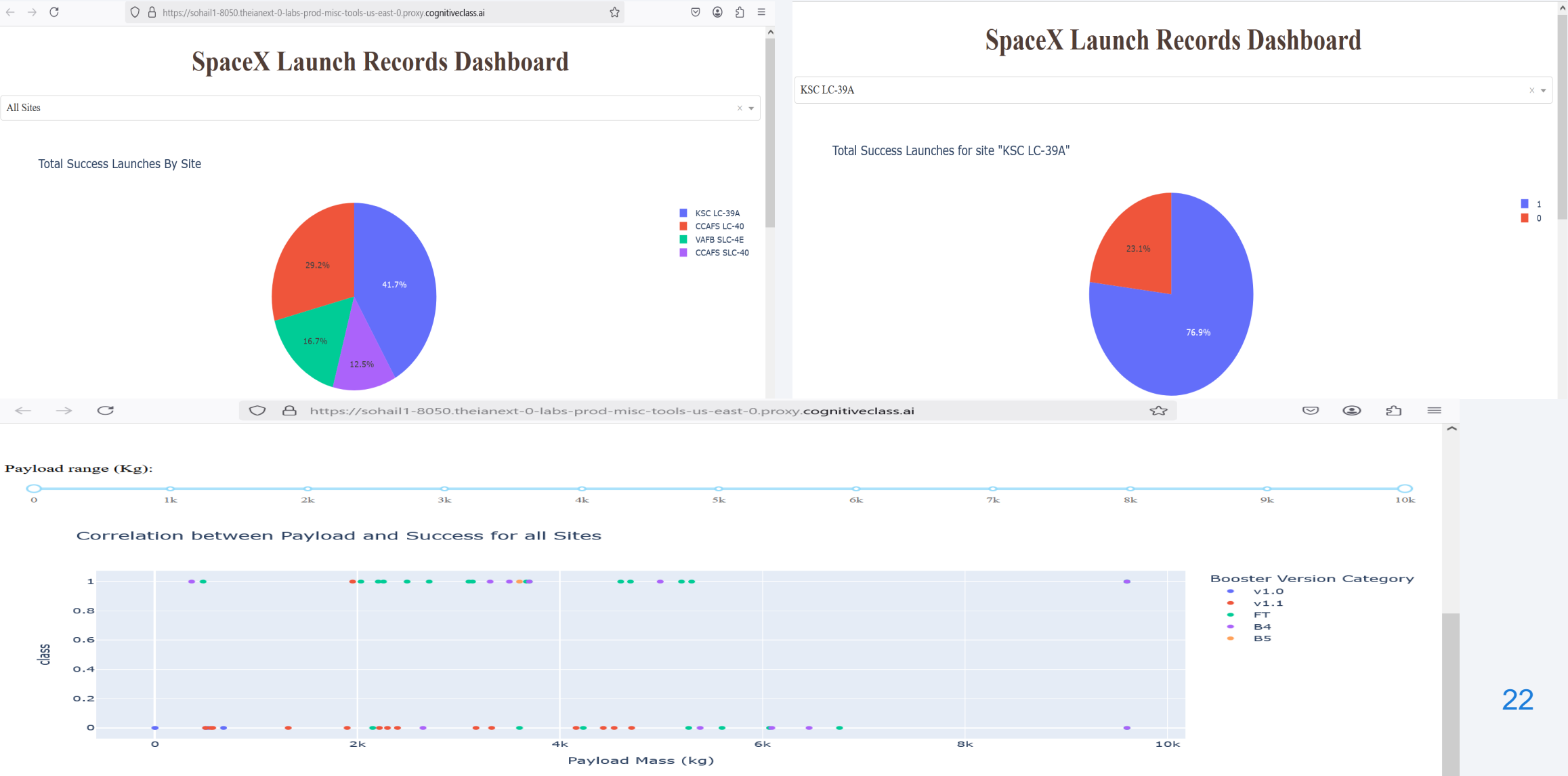
6. Selection of Best Performing Model:

- Identified the model with the highest accuracy on the test set as the best performer.
- Considered both training and test set performance to avoid overfitting and ensure real-world applicability.

Predictive Analysis (Flowchart)



Results

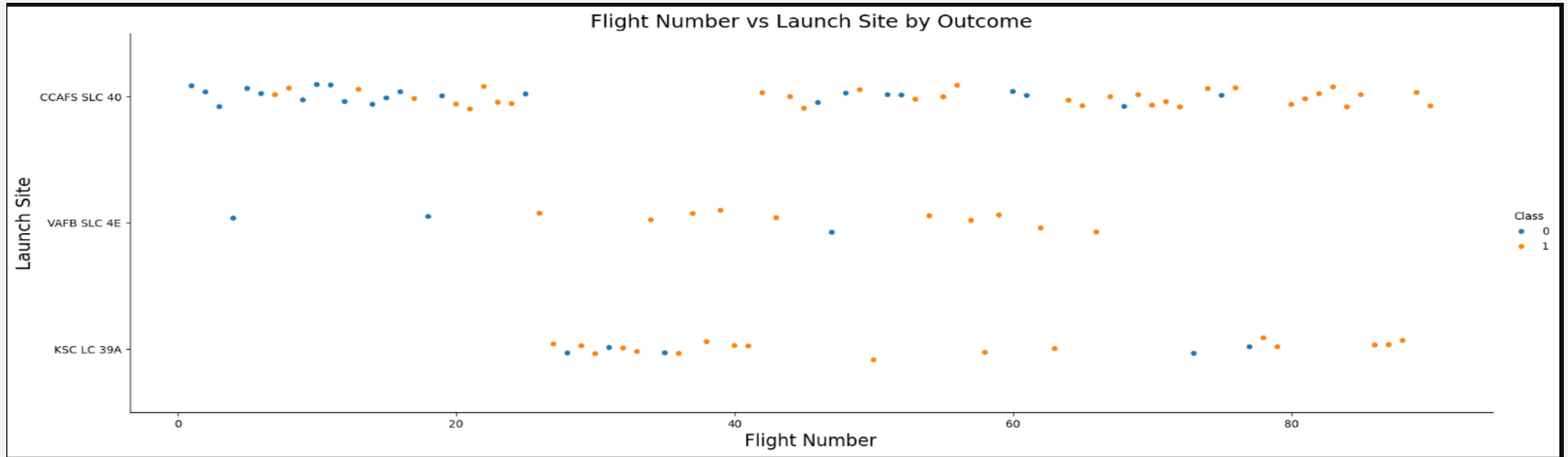


The background of the slide is an abstract composition. It features a dark blue base color. Overlaid on this are numerous diagonal streaks in shades of red and cyan. A faint, light blue grid pattern is also visible, particularly in the lower-left quadrant. The overall effect is dynamic and technological.

Section 2

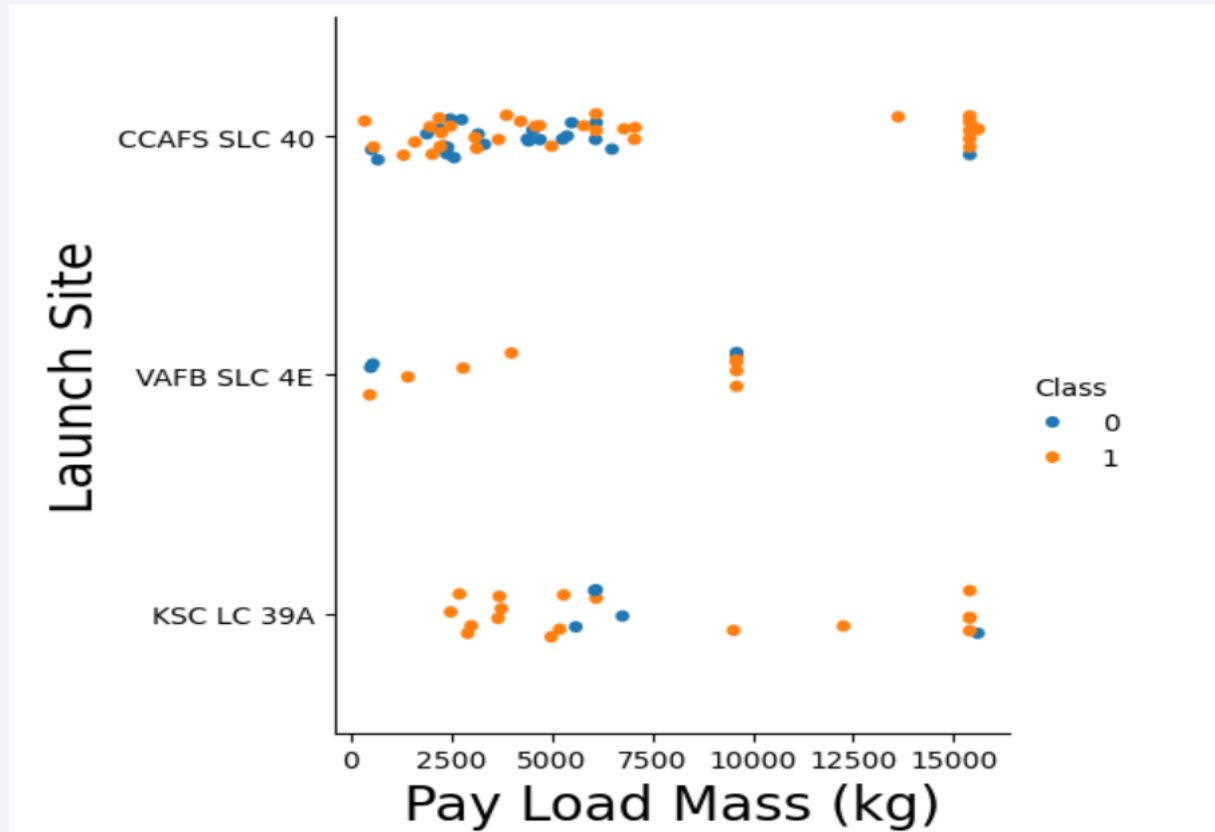
Insights drawn from EDA

Flight Number vs. Launch Site



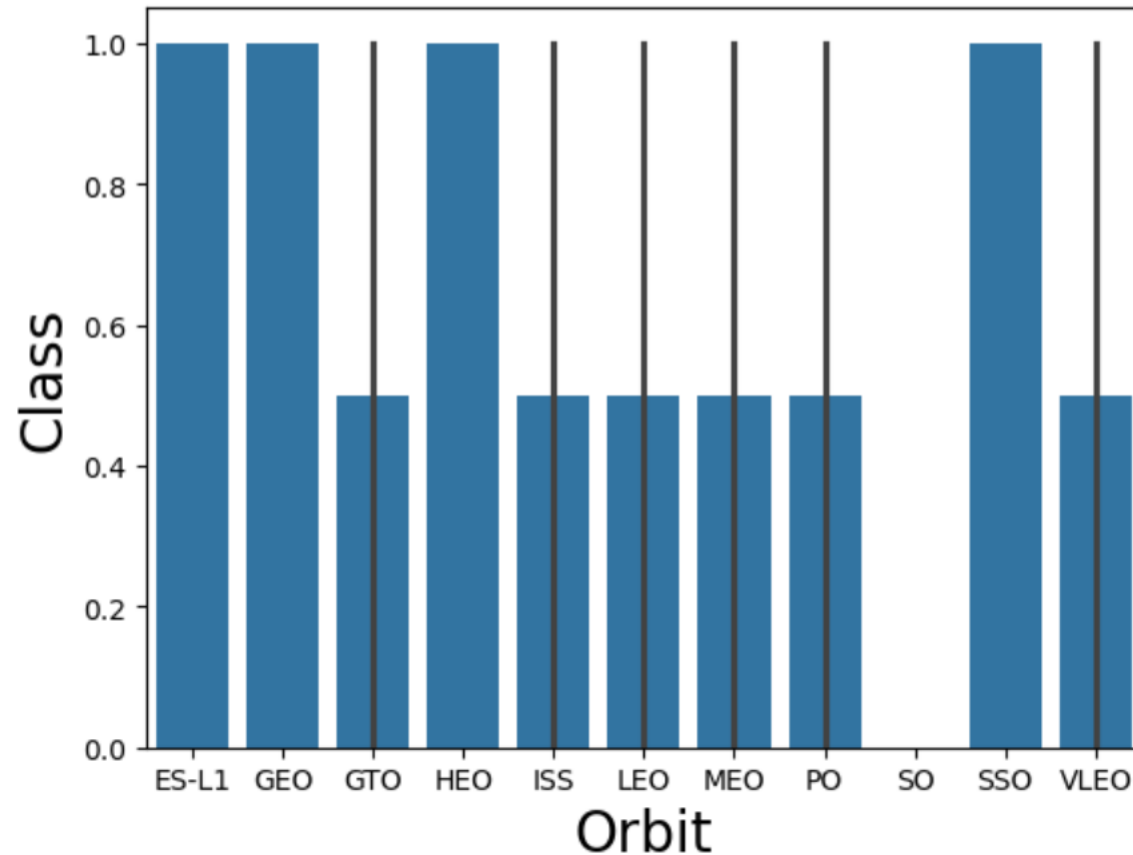
- **Mixed Outcomes at Major Launch Sites:** Both CCAFS SLC 40 and KSC LC 39A have a mix of successful (orange) and unsuccessful (blue) landings, indicating that factors other than the launch site itself may influence the landing success.
- **Consistent Activity Across Flight Numbers:** Launches are spread across a wide range of flight numbers at all sites, suggesting consistent activity over time without a clear trend of increasing or decreasing landing success.

Payload vs. Launch Site



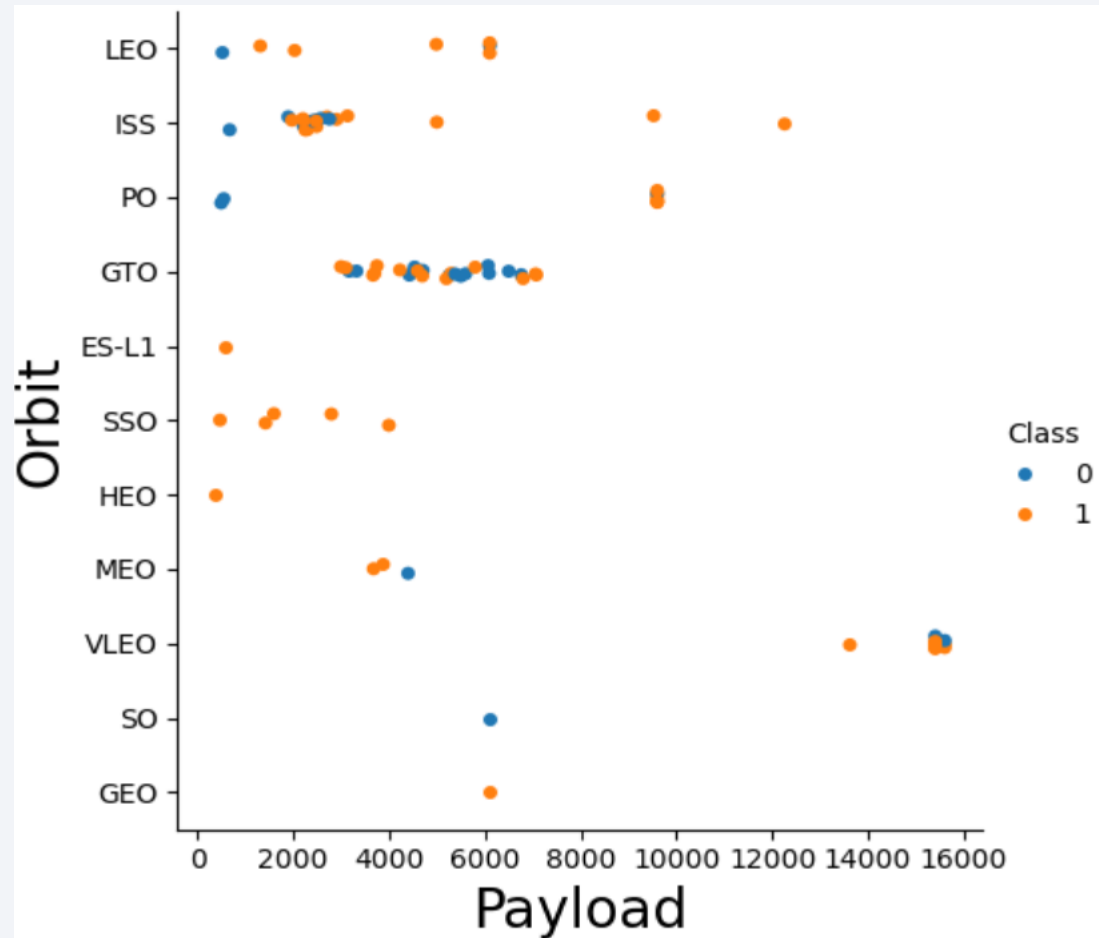
- **Payload Distribution:** Most launches from the CCAFS SLC 40 site handle payloads below 10,000 kg, while the VAFB SLC 4E and KSC LC 39A sites have a wider range of payload masses, indicating varied mission profiles.
- **High-Capacity Launches:** The KSC LC 39A site is frequently used for launching heavier payloads, with multiple launches carrying over 15,000 kg, suggesting its suitability for high-capacity missions.

Success Rate vs. Orbit Type



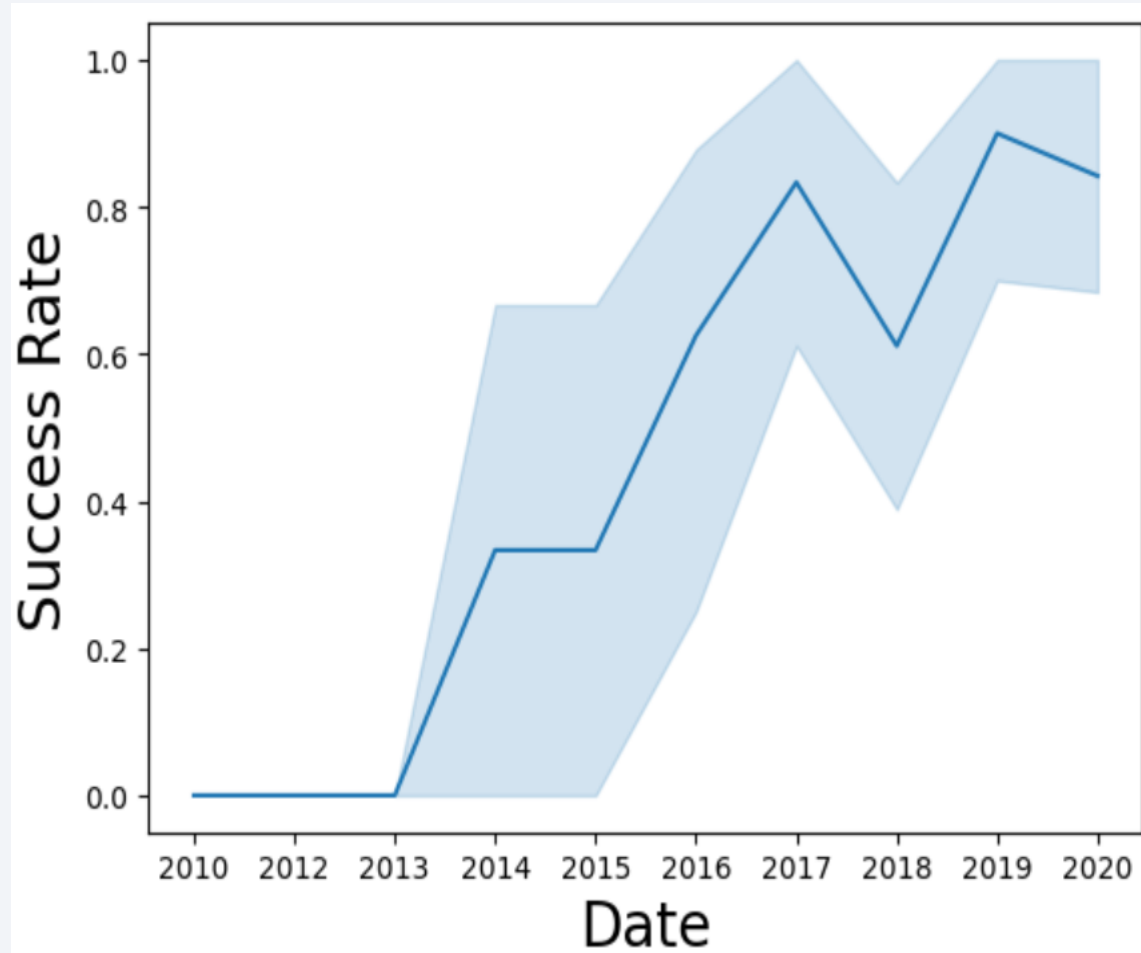
- **High Success Rates:** Missions to ES-L1, GEO, HEO, and SSO orbits have achieved a perfect success rate, indicating these orbits are highly reliable for successful first stage landings.
- **Lower Success Rate :** The GTO, ISS, LEO, MEO, PO, and VLEO orbit type shows a significantly lower success rate compared to other orbit types, suggesting that missions to this orbit may involve greater challenges or complexities.

Payload vs. Orbit Type



- **Increased Success Over Time:** The success rate of Falcon 9 launches improves significantly with higher flight numbers, indicating that experience and iterative improvements contribute to better outcomes.
- **Orbit-Specific Performance:** Early flights to GTO and ISS orbits had mixed outcomes, but recent missions to these orbits show a higher success rate, reflecting advancements in mission planning and execution.

Launch Success Yearly Trend



- The annual launch success rate has shown a significant improvement from 2013 onwards, reaching over 80% by 2020.
- Despite a dip in 2018, the overall trend indicates increasing reliability and success in Falcon 9 launches over the years.

All Launch Site Names

Task 1

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

```
%sql select distinct(Launch_Site) from SPACEXTABLE
```

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

Done.

Launch_Site

CCAFS LC-40

VAFB SLC-4E

KSC LC-39A

CCAFS SLC-40

Launch Site Names Begin with 'CCA'

Task 2

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

```
%sql select * from SPACEXTABLE where LAUNCH_SITE like 'CCA%' limit 5
```

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

Date	Time (UTC)	Booster_Version	Launch_Site	Payload	PAYLOAD_MASS_KG_	Orbit	Customer	Mission_Outcome	Landing_Outcome
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

Task 3

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

```
%sql select sum(PAYLOAD_MASS__KG_) from SPACEXTABLE where CUSTOMER = 'NASA (CRS)'
```

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

Done.

sum(PAYLOAD_MASS__KG_)

45596

Average Payload Mass by F9 v1.1

Task 4

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

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

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

Done.

avg(PAYLOAD_MASS__KG_)

2928.4

First Successful Ground Landing Date

Task 5

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

```
%sql select min(DATE) from SPACEXTABLE where Landing_Outcome = 'Success (ground pad)'
```

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

Done.

min(DATE)

2015-12-22

Successful Drone Ship Landing with Payload between 4000 and 6000

Task 6

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

```
%sql select BOOSTER_VERSION from SPACEXTBL where Landing_Outcome = 'Success (drone ship)' and PAYLOAD_MASS__KG_ > 4000 and PAYLOAD_MASS__KG_ < 6000
```

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

Done.

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

Total Number of Successful and Failure Mission Outcomes

Task 7

List the total number of successful and failure mission outcomes

```
%sql select count(MISSION_OUTCOME) from SPACEXTBL where MISSION_OUTCOME = 'Success' or MISSION_OUTCOME = 'Failure (in flight)'
```

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

```
Done.
```

```
count(MISSION_OUTCOME)
```

```
99
```

Boosters Carried Maximum Payload

Task 8

List the names of the booster_versions which have carried the maximum payload mass. Use a subquery

```
%sql select BOOSTER_VERSION from SPACEXTBL where PAYLOAD_MASS_KG_ = (select max(PAYLOAD_MASS_KG_) from SPACEXTBL)
```

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

Done .

Booster_Version

F9 B5 B1048.4

F9 B5 B1049.4

F9 B5 B1051.3

F9 B5 B1056.4

F9 B5 B1048.5

F9 B5 B1051.4

F9 B5 B1049.5

F9 B5 B1060.2

F9 B5 B1058.3

F9 B5 B1051.6

F9 B5 B1060.3

F9 B5 B1049.7

2015 Launch Records

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

```
%%sql
SELECT
  CASE
    WHEN substr("Date", 6, 2) = '01' THEN 'January'
    WHEN substr("Date", 6, 2) = '02' THEN 'February'
    WHEN substr("Date", 6, 2) = '03' THEN 'March'
    WHEN substr("Date", 6, 2) = '04' THEN 'April'
    WHEN substr("Date", 6, 2) = '05' THEN 'May'
    WHEN substr("Date", 6, 2) = '06' THEN 'June'
    WHEN substr("Date", 6, 2) = '07' THEN 'July'
    WHEN substr("Date", 6, 2) = '08' THEN 'August'
    WHEN substr("Date", 6, 2) = '09' THEN 'September'
    WHEN substr("Date", 6, 2) = '10' THEN 'October'
    WHEN substr("Date", 6, 2) = '11' THEN 'November'
    WHEN substr("Date", 6, 2) = '12' THEN 'December'
    ELSE 'Unknown'
  END AS "Month_Name",
  "Mission_Outcome",
  "Booster_Version",
  "Launch_Site"
FROM
  SPACEXTABLE
WHERE
  substr("Date", 0, 5) = '2015';
```

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

Month_Name	Mission_Outcome	Booster_Version	Launch_Site
January	Success	F9 v1.1 B1012	CCAFS LC-40
February	Success	F9 v1.1 B1013	CCAFS LC-40
March	Success	F9 v1.1 B1014	CCAFS LC-40
April	Success	F9 v1.1 B1015	CCAFS LC-40
April	Success	F9 v1.1 B1016	CCAFS LC-40
June	Failure (in flight)	F9 v1.1 B1018	CCAFS LC-40
December	Success	F9 FT B1019	CCAFS LC-40

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

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.

```
%%sql
```

```
SELECT
    "Landing_Outcome",
    COUNT(*) AS "Count"
FROM
    SPACEXTABLE
WHERE
    "Date" BETWEEN '2010-06-04' AND '2017-03-20'
GROUP BY
    "Landing_Outcome"
ORDER BY
    COUNT(*) DESC;
```

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

```
:  Landing_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
```

A satellite view of Earth from space, showing the curvature of the planet and city lights at night. The background is a deep blue gradient.

Section 3

Launch Sites Proximities Analysis

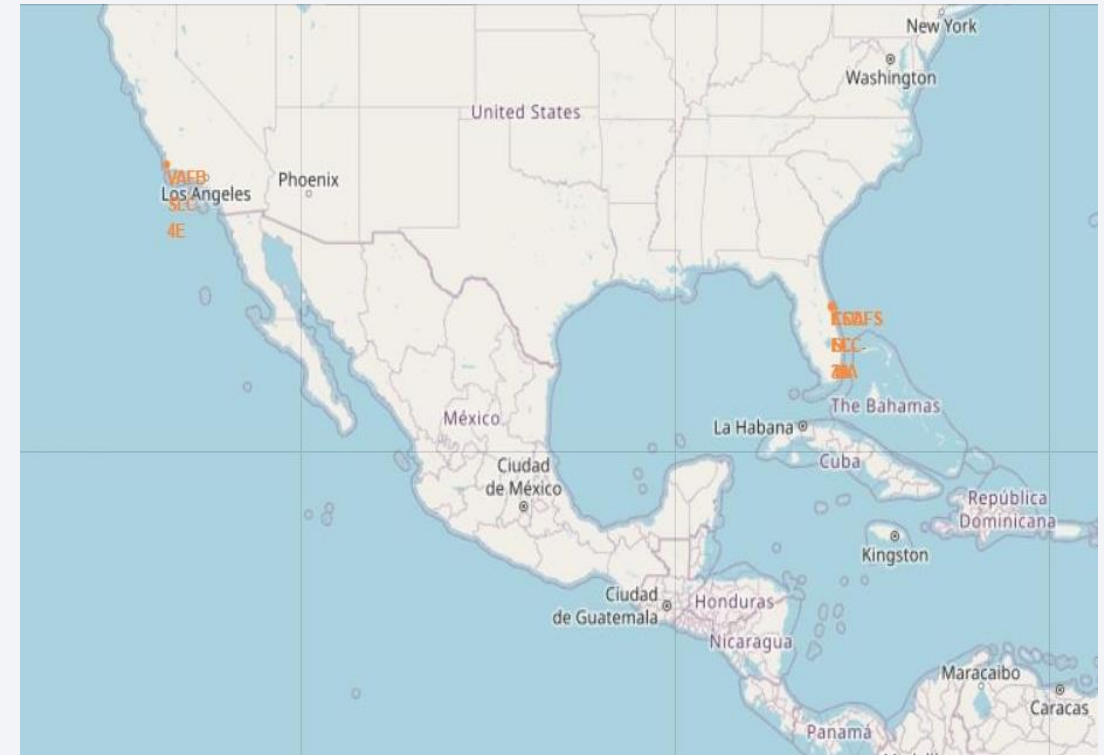
Task 1: Mark all launch sites on a map

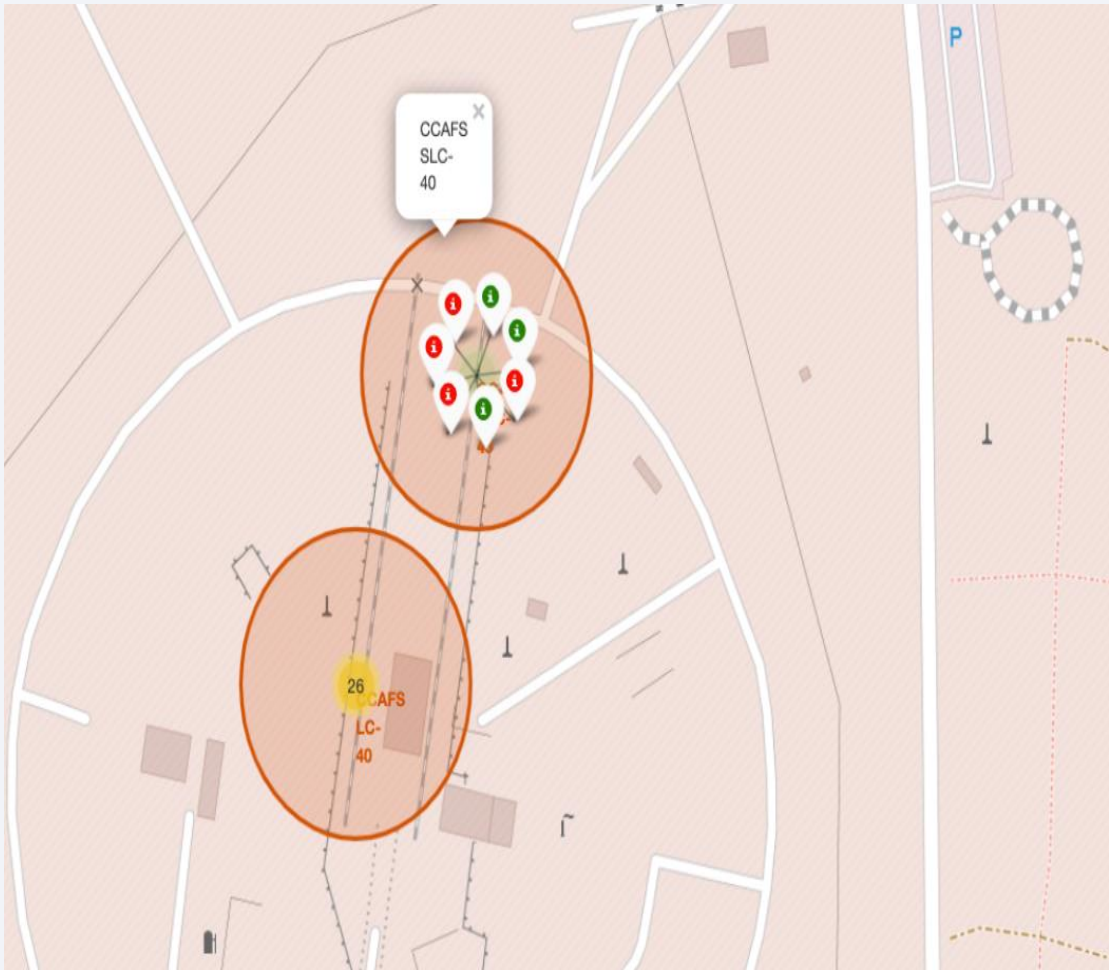
1. Are all launch sites in proximity to the Equator line?

- No, not all launch sites are in close proximity to the Equator.
- The launch site at Vandenberg Air Force Base (VAFB SLC-4E) is located at a latitude of 34.63, which is further from the Equator compared to the other sites in Florida.

2. Are all launch sites in very close proximity to the coast?

- Yes, all launch sites are in close proximity to the coast.
- The Cape Canaveral sites (CCAFS LC-40 and CCAFS SLC-40) and Kennedy Space Center (KSC LC-39A) are near the coast in Florida.
- Vandenberg Air Force Base (VAFB SLC-4E) is also near the coast in California.





Task 3: Calculate the distances between a launch site to its proximities

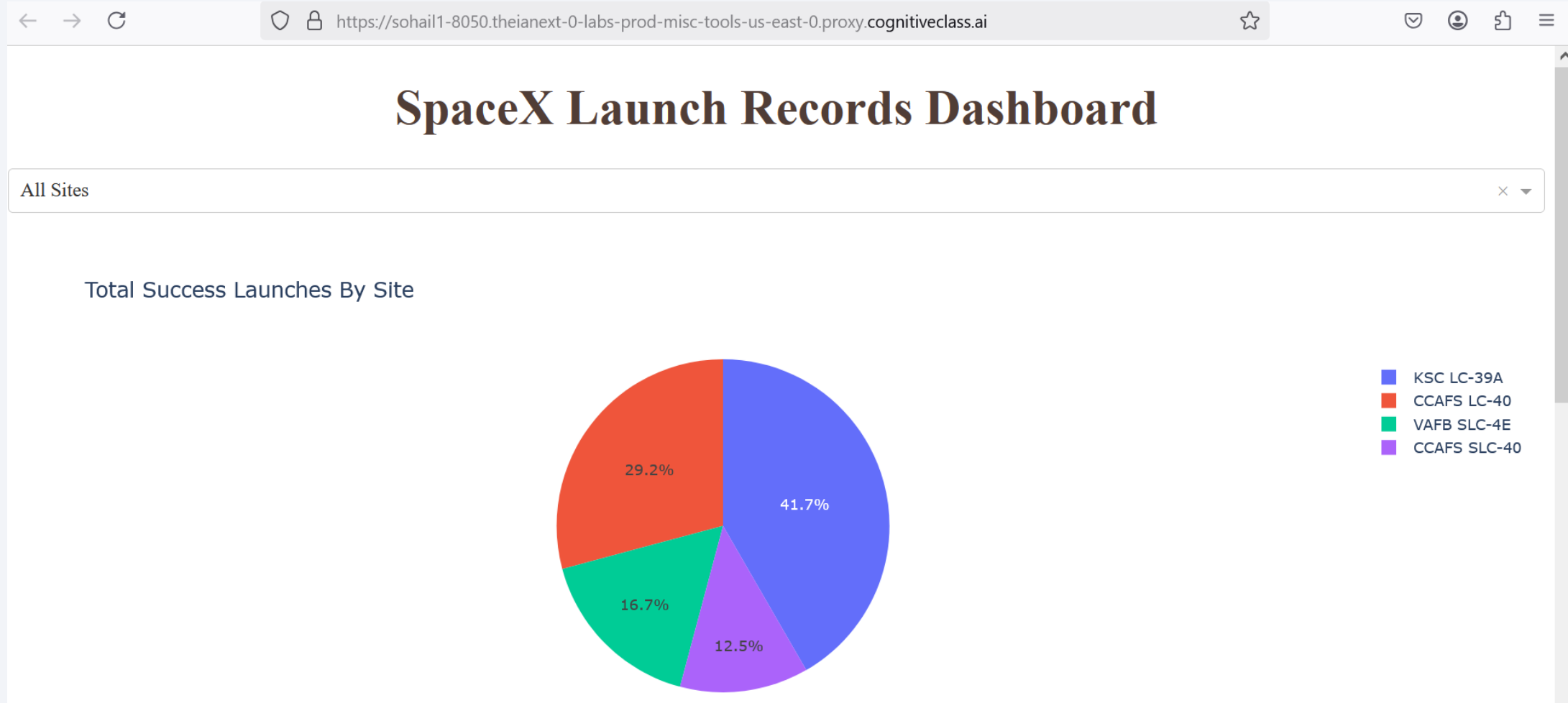




Section 4

Build a Dashboard with Plotly Dash

SpaceX Launch Records Dashboard



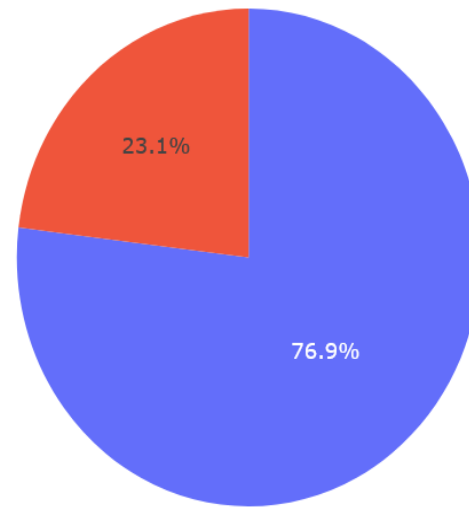
Pie chart for the launch site with highest launch success ratio

SpaceX Launch Records Dashboard

KSC LC-39A



Total Success Launches for site "KSC LC-39A"



■ 1
■ 0

Payload vs. Launch Outcome scatter plot for all sites

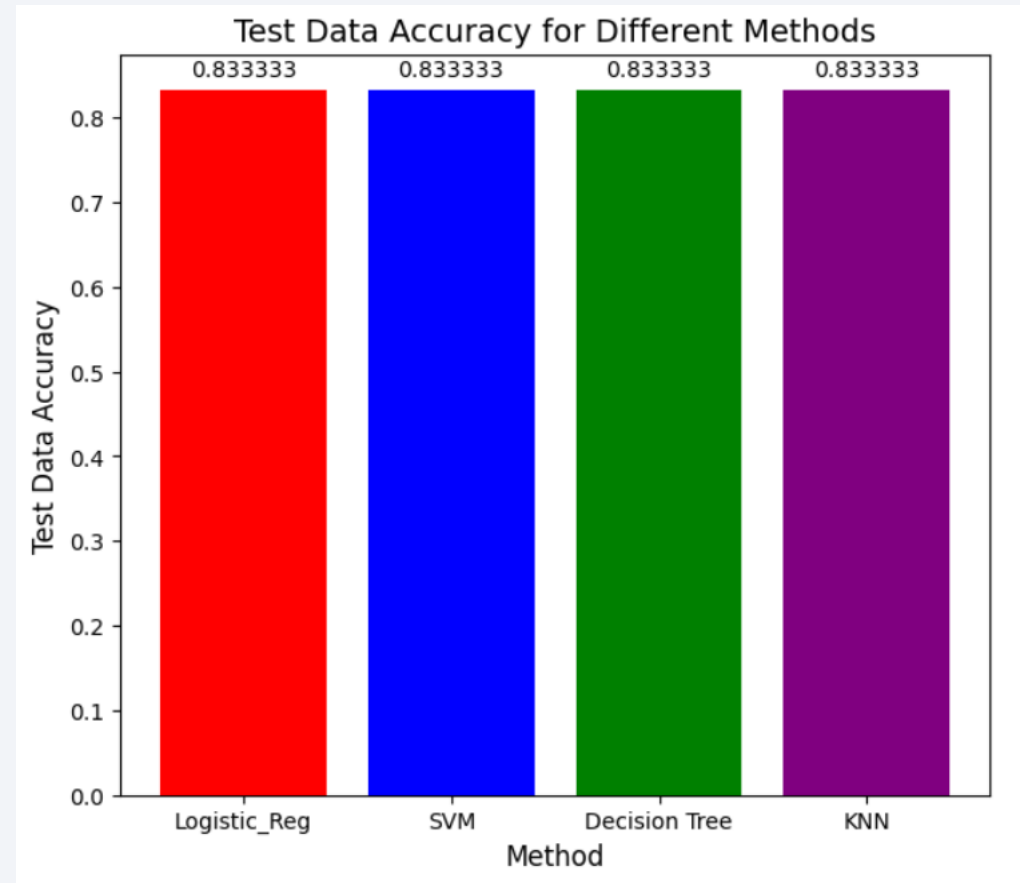


Section 5

Predictive Analysis (Classification)

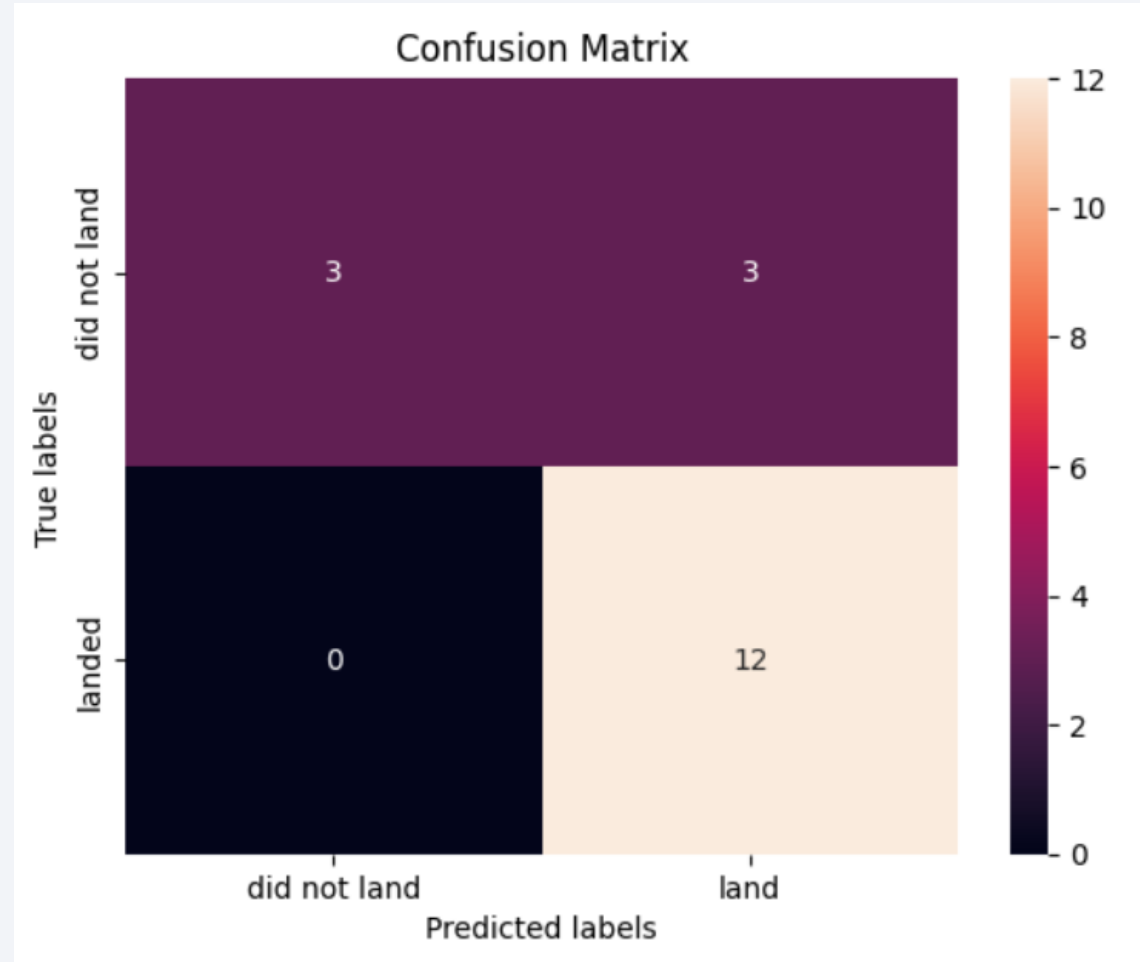
Classification Accuracy

- Based on the results, all the models like Logistic Regression, Support Vector Machine, Decision Tree and K Nearest Neighbors achieved the same accuracy of 0.8333.



Confusion Matrix

- Based on the results, all the models like Logistic Regression, Support Vector Machine, Decision Tree and K Nearest Neighbors achieved the same accuracy of 0.8333 and also have same confusion matrix.



Conclusions

Point 1: Our analysis revealed that the "KSC LC-39A" launch site has the highest success rate among all sites, accounting for 41.7% of successful launches. This indicates that this site might have optimal conditions or processes that contribute to a higher success rate.

Point 2: The scatter plot analysis showed that the "FT" booster version has a high success rate across various payload masses, demonstrating its reliability and robustness compared to other booster versions. This suggests that future missions might benefit from utilizing this booster version for improved success rates.

Point 3: No clear pattern was observed linking higher payload masses to lower success rates, indicating that factors other than payload mass, such as launch site conditions and booster versions, play a more significant role in determining the outcome of a launch.

Thank you!

