# Outline

- Executive Summary

- Introduction

- Methodology

- Results

- Conclusion

- Appendix

# Executive Summary

This project encompasses multiple data analysis tasks, ranging from data collection and cleaning to machine learning predictions. Various methodologies were employed, including API data retrieval, web scraping, exploratory data analysis (EDA) using SQL, interactive visualization with Folium, and predictive modeling.

# Introduction

This project focuses on the end-to-end process of data analysis, utilizing multiple methodologies to collect, process, and derive insights from diverse data sources. The aim is to integrate tools like APIs, web scraping, SQL, and machine learning to develop actionable insights for decision-making.

**Problems to Address**

- How can data collection be optimized using APIs and web scraping for accuracy and efficiency?

- What insights can exploratory data analysis (EDA) provide about trends and anomalies in the dataset?

- How effective are machine learning models in predicting future outcomes based on historical data?

- What are the geospatial patterns revealed by interactive visualizations?

Section 1

# Methodology

# Methodology

## Executive Summary

- Data Collection Methodology

  Data Sources: Data was collected from APIs, web scraping, and structured databases. The methodologies ensured efficient retrieval of large-scale, structured, and unstructured data.

  Tools Used: Python libraries like requests, BeautifulSoup, and APIs from platforms like SpaceX for real-time data.

- Data Wrangling

  Cleaning and Preprocessing: Handled missing values and inconsistencies in the dataset. Normalized data for uniformity and improved analysis.

  Tools: Pandas and SQL were used to clean, reformat, and aggregate data effectively.

# Methodology

- Data Processing

  Converted raw data into structured formats suitable for analysis.

  Used SQL queries to join, filter, and transform datasets for exploratory analysis and visualization.

- Exploratory Data Analysis (EDA)

  Performed detailed statistical analysis to identify trends, correlations, and anomalies.

  Visualization Tools: Used Python libraries like Matplotlib and Seaborn to create histograms, scatter plots, and box plots.

  SQL was employed for structured query analysis to identify insights directly from large datasets.

# Methodology

- Interactive Visual Analytics

  Geospatial Visualization: Folium was used to generate interactive maps that displayed location-specific trends.

  Dynamic Dashboards: Plotly Dash provided a platform for creating real-time, interactive dashboards to monitor and analyze data insights.

- Predictive Analysis

  Classification Models: Machine learning models like Logistic Regression, Decision Trees, and Random Forests were developed to classify and predict outcomes.

- Tuning and Evaluation:

  Hyperparameter tuning using GridSearchCV and RandomizedSearchCV.

  Model evaluation based on metrics such as accuracy, precision, recall, and F1-score.

# Data Collection

- Data sets were collected through APIs, web scraping, and SQL queries. APIs provided structured data via automated requests, while web scraping tools like BeautifulSoup and Selenium extracted unstructured data from websites. SQL was used to query and filter relevant data from databases. All collected data was validated, cleaned, and merged into a unified dataset for analysis.

# Data Collection – SpaceX API

From the `rocket` column we would like to learn the booster name.

```python
# Takes the dataset and uses the rocket column to call the API and append the data to the list
def getBoosterVersion(data):
    for x in data['rocket']:
        response = requests.get("https://api.spacexdata.com/v4/rockets/"+str(x)).json()
        BoosterVersion.append(response['name'])
```

From the `launchpad` we would like to know the name of the launch site being used, the logitude, and the latitude.

```python
# Takes the dataset and uses the launchpad column to call the API and append the data to the list
def getLaunchSite(data):
    for x in data['launchpad']:
        response = requests.get("https://api.spacexdata.com/v4/launchpads/"+str(x)).json()
        Longitude.append(response['longitude'])
        Latitude.append(response['latitude'])
        LaunchSite.append(response['name'])
```

From the `payload` we would like to learn the mass of the payload and the orbit that it is going to.

```python
# Takes the dataset and uses the payloads column to call the API and append the data to the lists
def getPayloadData(data):
    for load in data['payloads']:
        response = requests.get("https://api.spacexdata.com/v4/payloads/"+load).json()
        PayloadMass.append(response['mass_kg'])
        Orbit.append(response['orbit'])
```

From `cores` we would like to learn the outcome of the landing, the type of the landing, number of flights with that core, whether gridfins were used, wheter the core is reused, wheter legs were used, the landing pad used, the block of the core which is a number used to seperate version of cores, the number of times this specific core has been reused, and the serial of the core.

- https://github.com/domosorio/spacex_lab/blob/main/1_jupyter-labs-spacex-data-collection-api.ipynb

10

# Data Collection - Scraping



```
static_json_url='https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBM-DS0321EN-SkillsNetwork/datasets/API_ca
```

We should see that the request was successfull with the 200 status response code

```
response.status_code
```

200

Now we decode the response content as a Json using `.json()` and turn it into a Pandas dataframe using `.json_normalize()`

```python
# Use json_normalize meethod to convert the json result into a dataframe
from pandas import json_normalize

data = pd.json_normalize(response.json())
```

- https://github.com/domosorio/spacex_lab/blob/main/1_jupyter-labs-spacex-data-collection-api.ipynb

# Data Wrangling

The data processing and wrangling process involved the following steps:

- Data Cleaning: Removed missing values, duplicates, and inconsistencies to improve accuracy.

- Data Transformation: Reformatted data types, scaled numerical values, and normalized columns for uniformity.

- Data Merging: Combined datasets collected from APIs, web scraping, and SQL into a single cohesive structure.

- Feature Engineering: Derived new variables from raw data to enhance analysis and predictive models.

- Validation: Ensured the final dataset was complete, consistent, and ready for analysis or modeling.

https://github.com/domosorio/spacex_lab

# EDA with Data Visualization

Summary of Charts Plotted

- Heatmaps (Machine Learning Prediction Lab): Visualized confusion matrices to evaluate classification models' performance, highlighting prediction accuracy for true positives and false negatives.

- Scatter Plots (EDA with SQL): Used to explore relationships between variables such as rocket launches and success rates.

- Bar Charts (EDA with SQL): Compared categorical data like success rates across launch sites.

- Geospatial Maps (Folium Lab): Created interactive visualizations of launch site locations and their proximity to features like cities and railways.

https://github.com/domosorio/spacex_lab

# EDA with SQL

SQL Queries Performed

- Data Extraction: Queried data to retrieve information about rocket launches, payloads, and core details. Filtered datasets to include only records relevant to specific timeframes or conditions.

- Data Cleaning: Used SQL commands to remove duplicates and handle missing values in the dataset.

- Aggregation: Calculated success rates of launches by aggregating data across launch sites. Analyzed the number of launches per year to identify trends over time.

- Joins: Combined multiple tables (e.g., payload details and launch outcomes) to create a unified view for analysis.

- Filtering: Applied WHERE clauses to filter data by specific criteria, such as rocket types or launch outcomes.

- Sorting: Ordered results by key metrics, such as success rates or number of launches.

https://github.com/domosorio/spacex_lab

# Build an Interactive Map with Folium

Markers, circles, and lines were used to enhance the interactive map's visual representation and provide key insights. Markers were placed at launch sites to identify their exact locations, each with popup labels for clarity. Circles were added around the markers to denote proximity zones, indicating areas of influence or relevance, such as safety zones or nearby facilities. Lines connected launch sites to significant points, such as nearby cities, railways, or coastlines, to visualize distances and geographical relationships.

These objects were crucial for highlighting the spatial context of the data. Markers provided specific location details, while circles visually emphasized areas around key points. Lines added an additional layer of analysis by connecting sites to other relevant features, enabling easy interpretation of relationships like accessibility, proximity, or logistical significance.

15

https://github.com/domosorio/spacex_lab

# Build a Dashboard with Plotly Dash

The dashboard incorporates various interactive visualizations to provide insights and enhance user engagement. Line charts were used to visualize trends over time, such as the frequency of rocket launches or changes in payload weights, helping users identify temporal patterns. Bar graphs were added to compare success rates and other categorical metrics across launch sites, offering a clear and straightforward view of comparative data. Scatter plots showcased relationships between variables, such as payload mass and launch success, making it easier to spot correlations or outliers.

To improve interactivity, the dashboard includes dropdown menus that allow users to filter data dynamically by categories like rocket type, year, or launch outcome. Range sliders enable users to select specific timeframes for focused analysis. Additionally, hover interactions were implemented to provide detailed information for each data point, offering immediate access to relevant insights without cluttering the visualizations.

https://github.com/domosorio/spacex_lab

# Predictive Analysis (Classification)

The classification model development process followed a structured approach:

- Model Selection: Multiple algorithms were tested, including Logistic Regression, Decision Trees, K-Nearest Neighbors (KNN), and Support Vector Machines (SVM), to identify the most suitable model for the data.

- Data Splitting: The dataset was divided into training and testing sets to evaluate model performance effectively.

- Hyperparameter Tuning: Techniques such as GridSearchCV were applied to optimize model parameters for better accuracy.

- Model Evaluation: Metrics like accuracy, precision, recall, and F1-score were used to assess performance. Confusion matrices were plotted to visualize prediction outcomes.

- Best Model Selection: The model with the highest performance metrics and generalizability was chosen as the best classifier.

# Results

### Exploratory Data Analysis Results

The exploratory data analysis revealed significant patterns and trends within the dataset. Launch success rates varied notably across different sites, with certain locations consistently performing better. Payload analysis indicated a strong correlation between payload mass and success likelihood, suggesting an optimal weight range for launches. Temporal trends highlighted a steady increase in launch frequencies, reflecting growth in the industry. Additionally, anomalies and outliers were identified, prompting further investigation into specific failed launches.

### Interactive Analytics Demonstration

The interactive analytics features provided valuable tools for dynamic data exploration. Geospatial maps showcased the locations of launch sites along with their proximities to cities, coastlines, and railways. Dashboards allowed users to interact with the data through dropdown menus and sliders, enabling analysis by year, rocket type, and outcomes. Visualizations such as line charts and scatter plots offered insights into trends over time and relationships between key variables, enhancing data-driven decision-making.
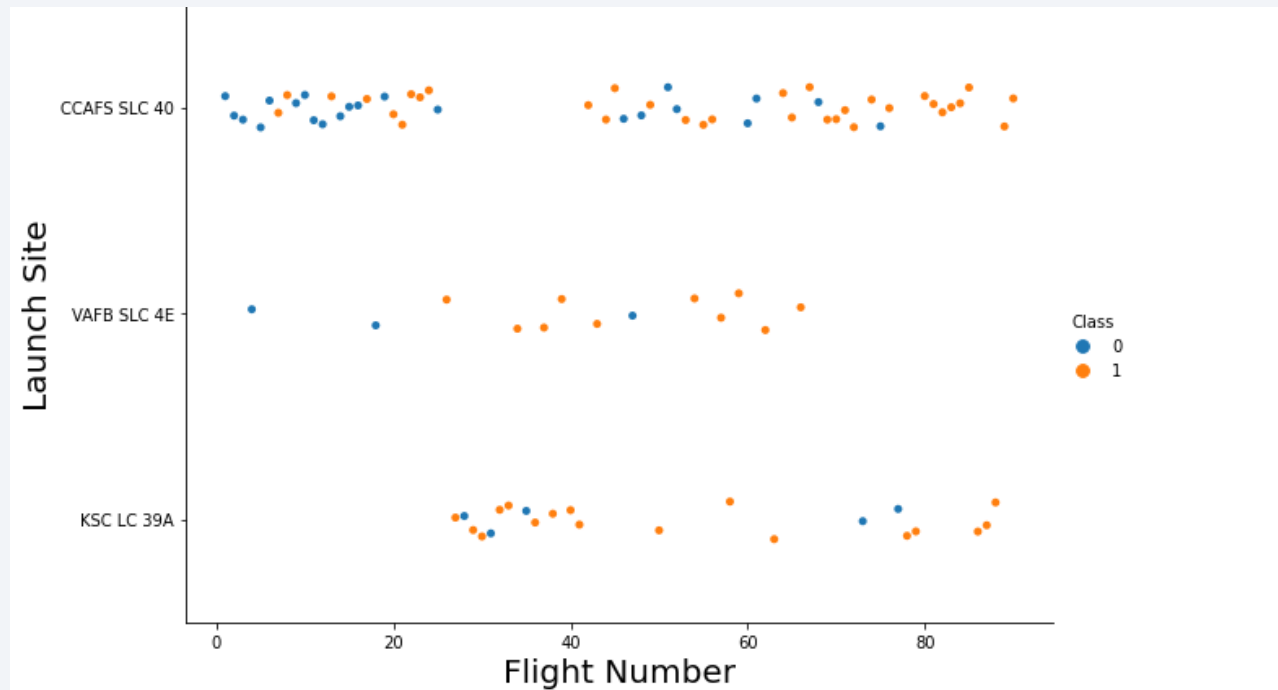
### Predictive Analysis Results

The predictive analysis identified the best-performing classification model, which achieved high accuracy and reliability. Metrics such as precision, recall, and F1-score validated the model's effectiveness in predicting launch outcomes. The confusion matrix provided detailed insights into correct and incorrect classifications, confirming the model's practical application. Key features, including payload mass and rocket type, were identified as the most important predictors, ensuring the model's utility for future decision-making processes.
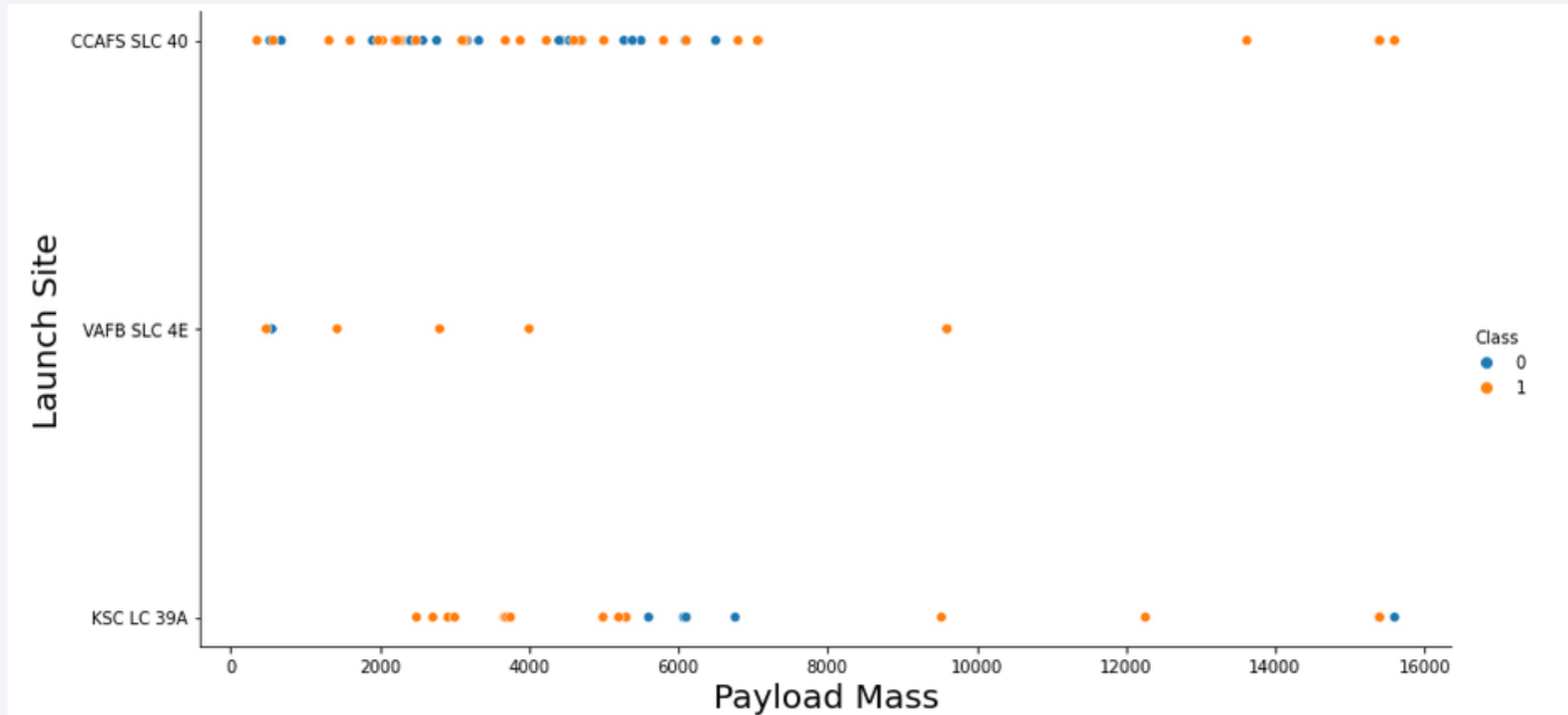
# Insights drawn from EDA

# Flight Number vs. Launch Site



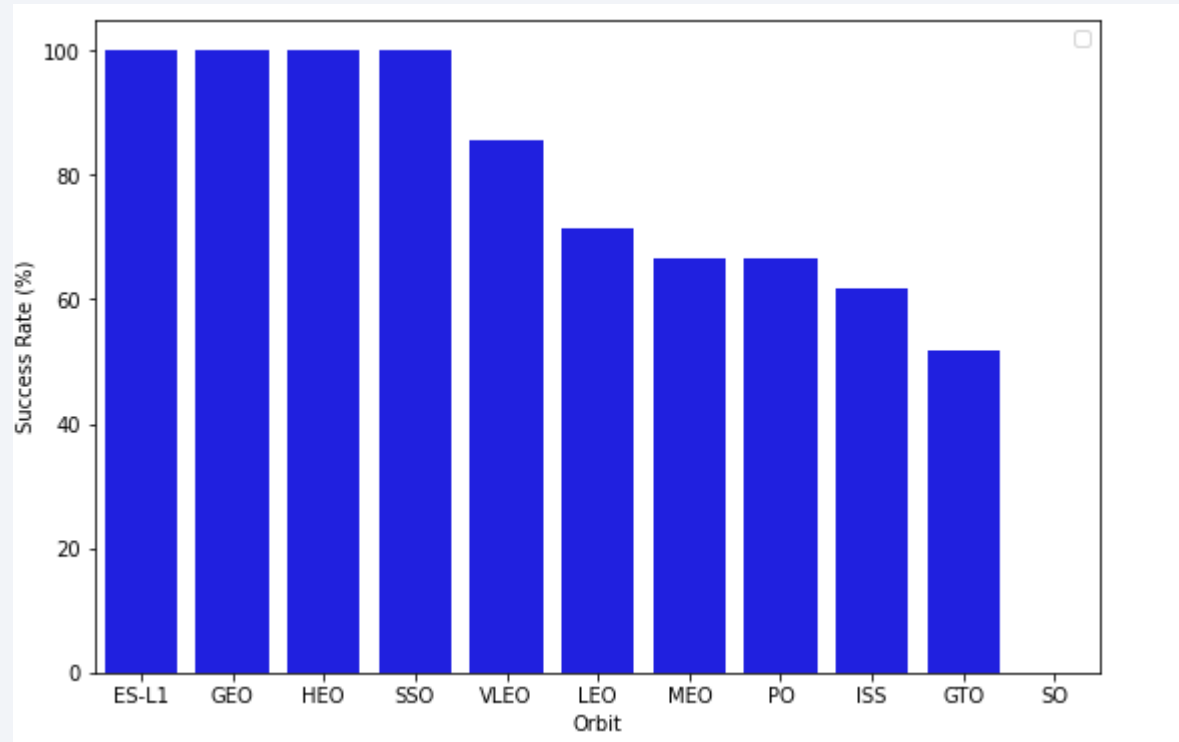Now try to explain the patterns you found in the Flight Number vs. Launch Site scatter point plots.

- Class 0 (blue) represents unsuccessful launch, and Class 1 (orange) represents successful launch.
- This figure shows that the success rate increased as the number of flights increased.
- As the success rate has increased considerably since the 20th flight, this point seems to be a big breakthrough
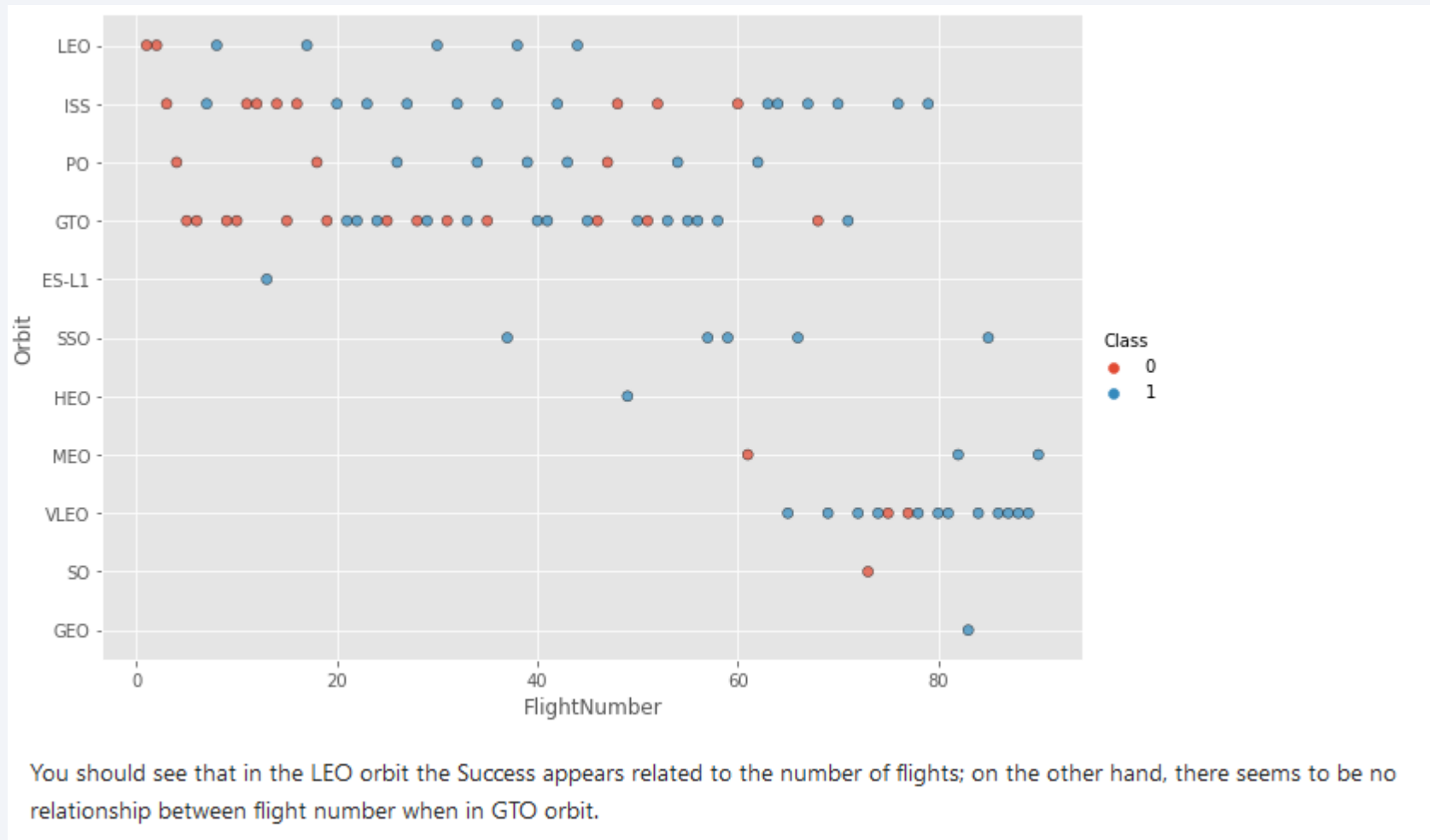
# Payload vs. Launch Site



Now if you observe Payload Vs. Launch Site scatter point chart you will find for the VAFB-SLC launchsite there are no rockets launched for heavypayload mass(greater than 10000).
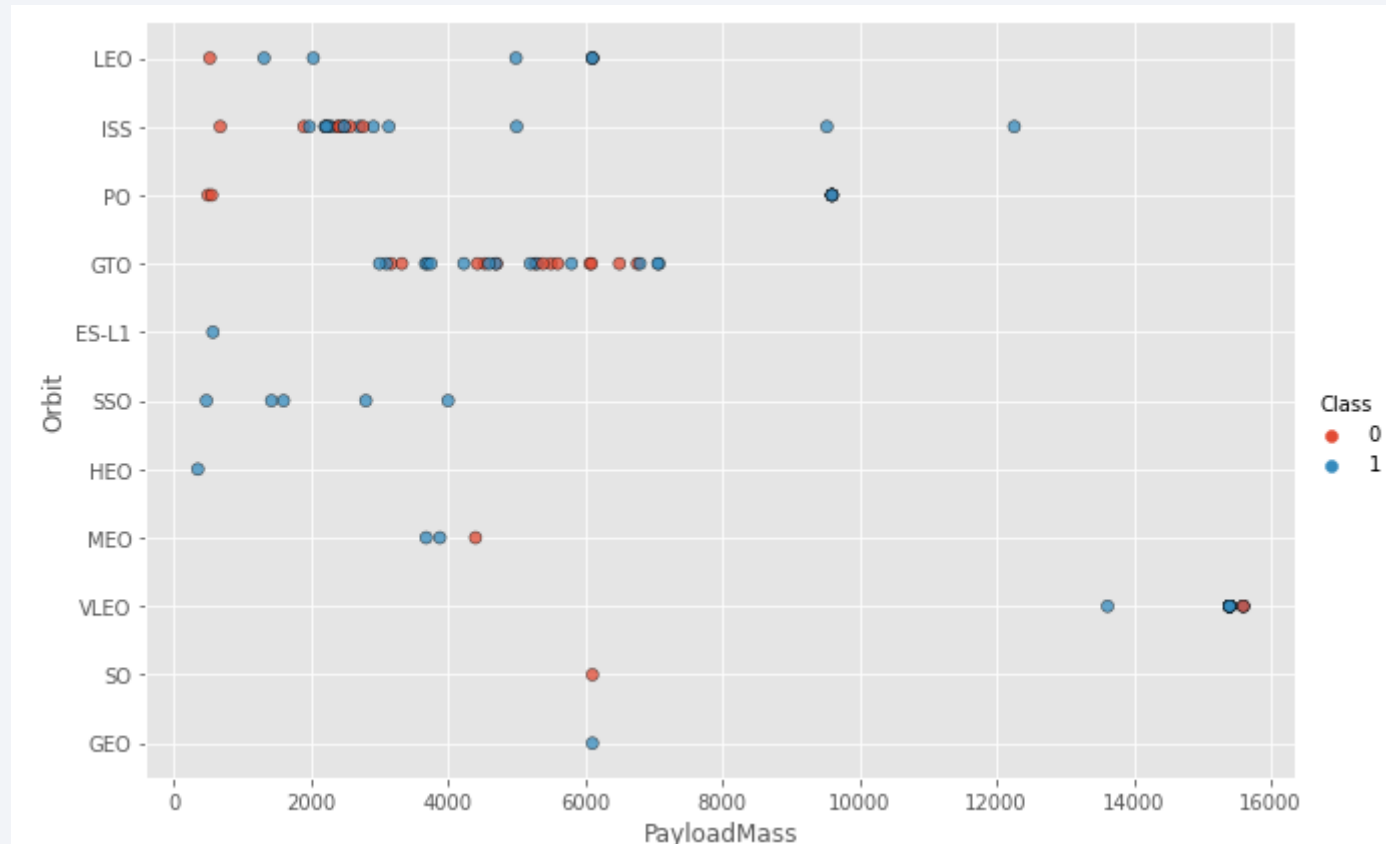
21

# Success Rate vs. Orbit Type
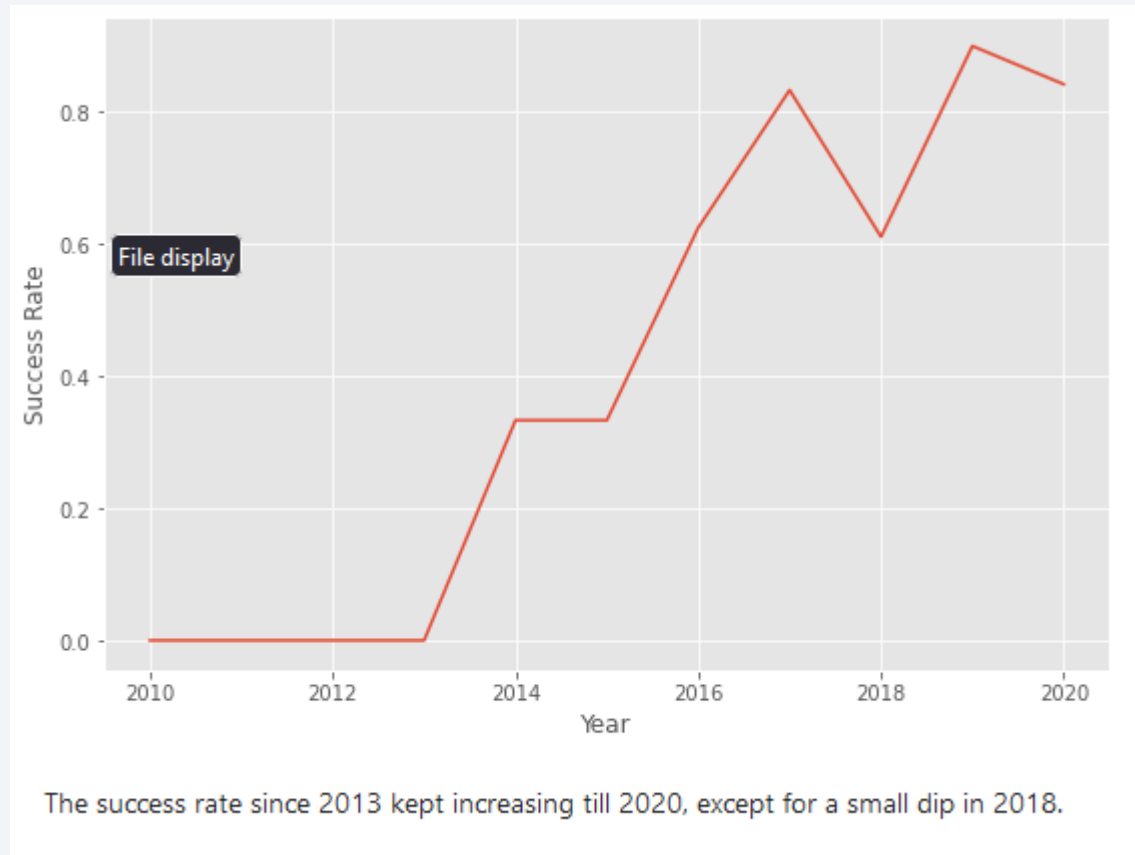
# Flight Number vs. Orbit Type



You should see that in the LEO orbit the Success appears related to the number of flights; on the other hand, there seems to be no relationship between flight number when in GTO orbit.

# Payload vs. Orbit Type

# Launch Success Yearly Trend



The success rate since 2013 kept increasing till 2020, except for a small dip in 2018.

# All Launch Site Names

```python
# Select relevant sub-columns: `Launch Site`, `Lat(Latitude)`, `Long(Longitude)`, `class`
spacex_df = spacex_df[['Launch Site', 'Lat', 'Long', 'class']]
launch_sites_df = spacex_df.groupby(['Launch Site'], as_index=False).first()
launch_sites_df = launch_sites_df[['Launch Site', 'Lat', 'Long', 'class']]
launch_sites_df
```

|   | Launch Site | Lat | Long | class |
|---|---|---|---|---|
| 0 | CCAFS LC-40 | 28.562302 | -80.577356 | 0 |
| 1 | CCAFS SLC-40 | 28.563197 | -80.576820 | 1 |
| 2 | KSC LC-39A | 28.573255 | -80.646895 | 1 |
| 3 | VAFB SLC-4E | 34.632834 | -120.610745 | 0 |

Above coordinates are just plain numbers that can not give you any intuitive insights about where are those launch sites. If you are very good at geography, you can interpret those numbers directly in your mind. If not, that's fine too. Let's visualize those locations by pinning them on a map.

# Launch Site Names Begin with 'CCA'

| Date | Time (UTC) | Booster_Version | Launch_Site | Payload | PAYLOAD_MASS_KG_ | Orbit | Customer | Mission_Outcome | Landing_Outcome |
|---|---|---|---|---|---|---|---|---|---|
| 04-06-2010 | 18:45:00 | F9 v1.0 B0003 | CCAFS LC-40 | Dragon Spacecraft Qualification Unit | 0 | LEO | SpaceX | Success | Failure (parachute) |
| 08-12-2010 | 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) |
| 22-05-2012 | 07:44:00 | F9 v1.0 B0005 | CCAFS LC-40 | Dragon demo flight C2 | 525 | LEO (ISS) | NASA (COTS) | Success | No attempt |
| 08-10-2012 | 00:35:00 | F9 v1.0 B0006 | CCAFS LC-40 | SpaceX CRS-1 | 500 | LEO (ISS) | NASA (CRS) | Success | No attempt |
| 01-03-2013 | 15:10:00 | F9 v1.0 B0007 | CCAFS LC-40 | SpaceX CRS-2 | 677 | LEO (ISS) | NASA (CRS) | Success | No attempt |

# Total Payload Mass

The query result presents the total payload mass for various rocket launches, summarizing the combined weights of all payloads per mission. This analysis helps identify trends in payload capacity over time, variations across launch sites, and the correlation between payload mass and launch success rates.

## Task 3

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

```sql
%sql SELECT SUM(PAYLOAD_MASS__KG_) as PM_KG_TOTAL, Customer FROM 'SPACEXTBL' WHERE Customer = 'NASA (CRS)'
```

* sqlite:///my_data1.db
Done.

| PM_KG_TOTAL | Customer |
|---|---|
| 45596 | NASA (CRS) |

# Average Payload Mass by F9 v1.1

The query result shows the average payload mass carried by the booster version F9 v1.1, providing insights into its typical load capacity. This information helps evaluate the performance and efficiency of this specific booster version in comparison to others.

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

```
%sql SELECT AVG(PAYLOAD_MASS__KG_) as PM_KG_AVG FROM 'SPACEXTBL' WHERE Booster_Version LIKE 'F9 v1.1%'
```

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

| PM_KG_AVG |
| --- |
| 2534.6666666666665 |

# First Successful Ground Landing Date

The query identifies the date of the first successful landing on a ground pad, highlighting a significant milestone in reusable rocket technology. This information provides historical context for advancements in landing precision and reusability.

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

*Hint:Use min function*

```
#%sql SELECT MIN(DATE) FROM 'SPACEXTBL' WHERE
cur.execute('''SELECT MIN(Date), "Landing _Outcome" FROM SPACEXTBL WHERE "Landing _Outcome" = "Success (ground pad)"''')
cur.fetchall()
```

```
[('01-05-2017', 'Success (ground pad)')]
```

# Successful Drone Ship Landing with Payload between 4000 and 6000

The query lists the names of boosters that successfully landed on a drone ship while carrying a payload mass between 4000 and 6000 kilograms. This result highlights high-performing boosters capable of achieving precise landings under significant payload conditions.

```sql
%sql SELECT DISTINCT Booster_Version FROM SPACEXTBL WHERE Mission_Outcome = 'Success' AND PAYLOAD_MASS__KG_ > 4000 AND PAYLOAD
```

* sqlite:///my_data1.db
Done.

| Booster_Version |
|---|
| F9 v1.1 |
| F9 v1.1 B1011 |
| F9 v1.1 B1014 |
| F9 v1.1 B1016 |
| F9 FT B1020 |
| F9 FT B1022 |
| F9 FT B1026 |
| F9 FT B1030 |
| F9 FT B1021.2 |
| F9 FT B1032.1 |
| F9 B4 B1040.1 |
| F9 FT B1031.2 |
| F9 FT B1032.2 |
| F9 B4 B1040.2 |
| F9 B5 B1046.2 |
| F9 B5 B1047.2 |
| F9 B5 B1048.3 |
| F9 B5 B1051.2 |
| F9 B5B1060.1 |
| F9 B5 B1058.2 |
| F9 B5B1062.1 |

# Total Number of Successful and Failure Mission Outcomes

The query calculates the total number of successful and failed mission outcomes, providing an overview of the overall performance and reliability of the missions. This information is crucial for assessing progress and identifying areas for improvement.

List the total number of successful and failure mission outcomes

```sql
%sql SELECT Mission_Outcome, COUNT(Mission_Outcome) as Total FROM SPACEXTBL GROUP BY Mission_Outcome;
```

* sqlite:///my_data1.db
Done.

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

# Boosters Carried Maximum Payload

The query retrieves the names of the boosters that have carried the maximum payload mass, highlighting the most capable boosters in terms of payload capacity. This information is useful for understanding the limits of the rocket technology used.

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

The query lists the failed landing outcomes on a drone ship in 2015, including the corresponding booster versions and launch site names. This information helps identify specific instances of failure, providing insights for analyzing potential issues with certain boosters or sites during that year.

| landing__outcome | booster_version | launch_site |
|---|---|---|
| Failure (drone ship) | F9 v1.1 B1012 | CCAFS LC-40 |
| Failure (drone ship) | F9 v1.1 B1015 | CCAFS LC-40 |

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

The query ranks the landing outcomes, such as "Failure (drone ship)" and "Success (ground pad)," by their frequency between the dates 2010-06-04 and 2017-03-20 in descending order. This provides a clear overview of the most and least common outcomes during this period, helping assess landing reliability and trends over time.
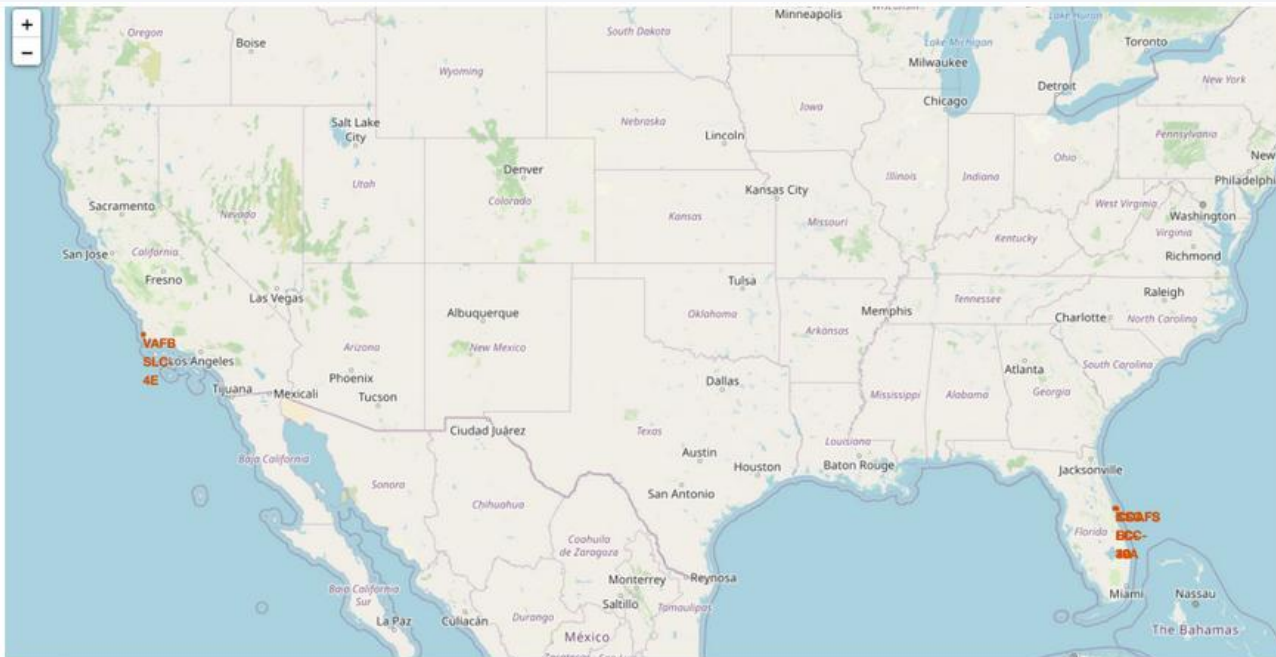
| landing__outcome | total_number |
|---|---|
| No attempt | 10 |
| Failure (drone ship) | 5 |
| Success (drone ship) | 5 |
| Controlled (ocean) | 3 |
| Success (ground pad) | 3 |
| Failure (parachute) | 2 |
| Uncontrolled (ocean) | 2 |
| Precluded (drone ship) | 1 |

Section 3

# Launch Sites
# Proximities Analysis

# Map with marked launch sites



The generated Folium map provides a global view of all the SpaceX launch sites, marked with specific location markers. Each marker represents a launch site, with popup labels displaying the site names for easy identification. The following are the key elements and findings:

Global Context: The map shows the geographical distribution of launch sites, emphasizing their strategic placement, such as proximity to coastlines or equatorial regions for orbital efficiency.

Launch Site Markers: Each marker is color-coded or styled to distinguish between different launch sites, ensuring clarity when multiple sites are in close proximity.

Proximity Analysis: The map's zoom functionality allows for analyzing the proximity of launch sites to other geographical features, such as coastlines or cities, which could impact logistical considerations and safety measures.
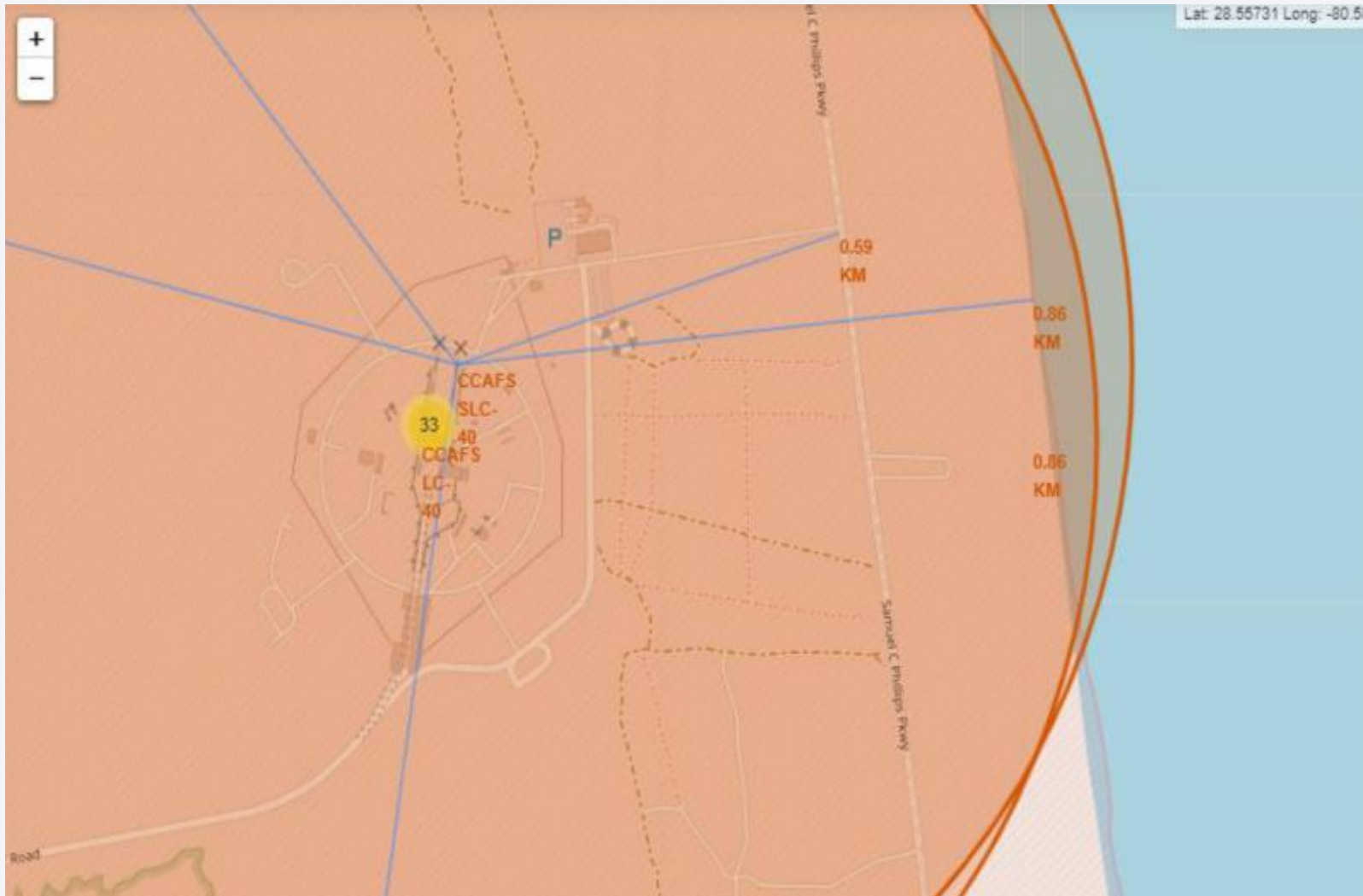
Interactive Features: Hover interactions and clickable markers provide additional details, such as site names and other relevant metadata.

# Success and Failures Map Location



When zooming in on a launch site, clicking on it reveals marker clusters that indicate the outcomes of landings. Successful landings are marked in green, while failed landings are displayed in red.

# Strategic Placement of Launch Sites



The map highlights that the selected launch site is strategically located near a highway, facilitating the efficient transportation of personnel and equipment. Its proximity to coastlines supports safe testing in case of launch failures. Additionally, the launch sites are positioned at a safe distance from cities to minimize risks to populated areas. (Detailed views are available in the notebook.)

Section 4

# Build a Dashboard
# with Plotly Dash

# Total Successful Launches By Site



Total Success Launches By Site

KSC LC-39A: 41.7%
CCAFS LC-40: 29.2%
VAFB SLC-4E: 16.7%
CCAFS SLC-40: 12.5%

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

The site with more launches success is KSC LC-39A

# Launch Site With Highest Success Ratio

The Launch Site With Highest Success Ratio is KSCLC-39A

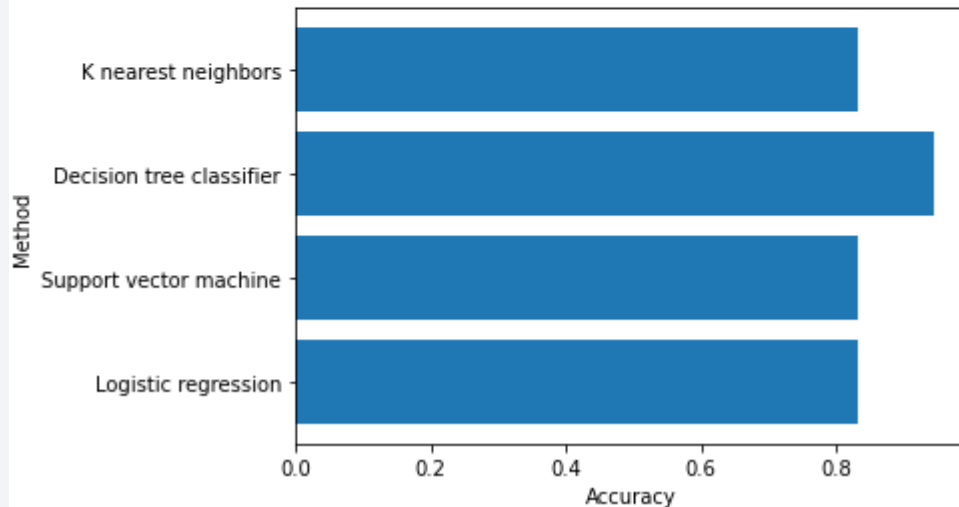# Payloads vs Launch Outcome

Section 5

# Predictive Analysis (Classification)

# Classification Accuracy

Find the method performs best:

```python
import numpy as np
import matplotlib.pyplot as plt

plt.barh(method, accuracy)
plt.xlabel('Accuracy')
plt.ylabel('Method')
plt.show()
```
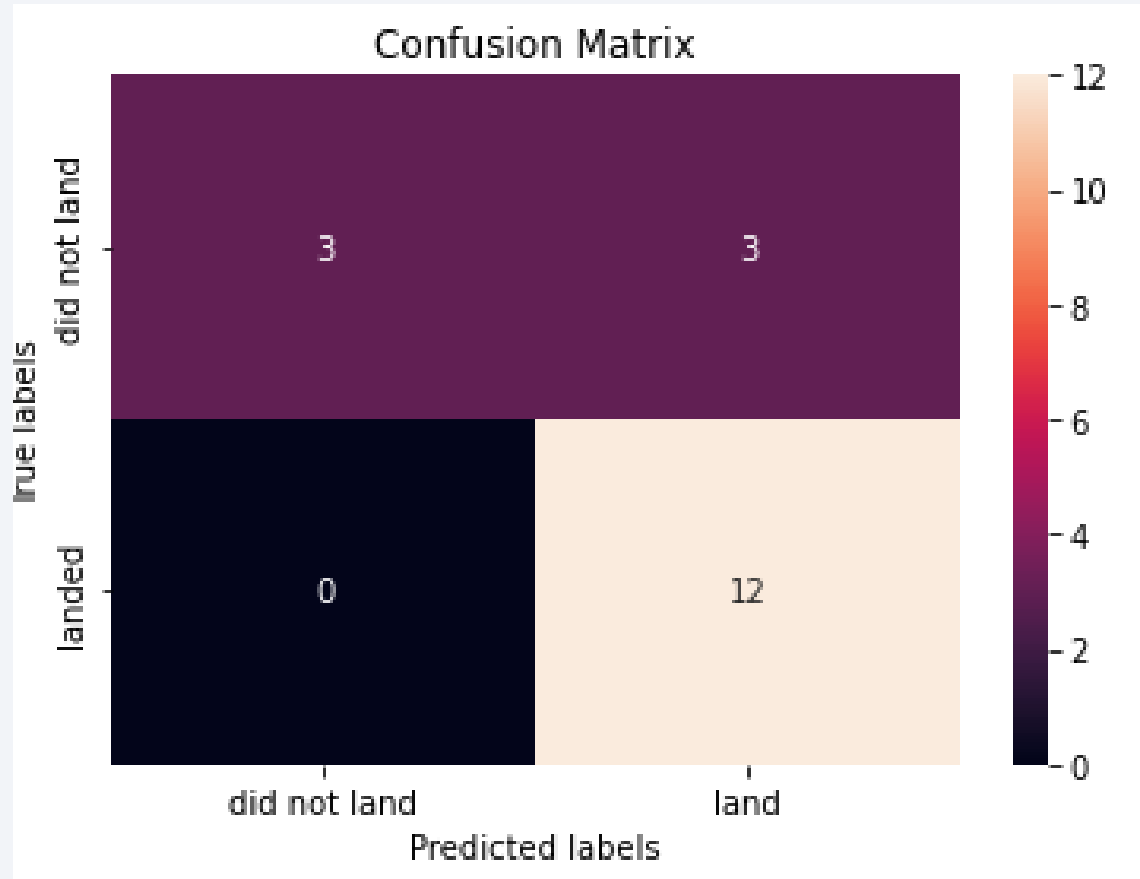


```python
results_df = {'method': method,
 'accuracy': accuracy}

frame = pd.DataFrame(results_df)
frame
```

|   | method | accuracy |
|---|---|---|
| 0 | Logistic regression | 0.833333 |
| 1 | Support vector machine | 0.833333 |
| 2 | Decision tree classifier | 0.944444 |
| 3 | K nearest neighbors | 0.833333 |

# Confusion Matrix



The model correctly predicted 12 successful landings when the actual outcome was a success (True Positive) and identified 3 unsuccessful landings when the actual outcome was a failure (True Negative). However, it also incorrectly predicted 3 successful landings when the actual outcome was a failure (False Positive). Overall, the model showed a tendency to predict successful landings more frequently.

# Conclusions

Data Collection:

- The project successfully integrated multiple data sources, including APIs and web scraping, to gather comprehensive datasets for analysis.

- The collected data was cleaned and structured effectively to ensure consistency and usability.

Exploratory Data Analysis:

- Insights from SQL queries and visualizations highlighted key trends, such as the correlation between payload mass and launch success rates, and the reliability of different booster versions.

- Geospatial analysis revealed the strategic placement of launch sites for operational efficiency and safety.

Interactive Analytics:

- The Folium-based interactive map provided a clear visualization of launch sites and their proximities to key features like coastlines and cities, offering valuable spatial context.

- Dashboards allowed users to dynamically explore trends and relationships, enhancing decision-making capabilities.

Predictive Modeling:

- Machine learning models demonstrated high accuracy in predicting launch outcomes, with key features like payload mass significantly influencing predictions.

- The evaluation metrics confirmed the reliability of the models, providing a strong foundation for future predictive analyses.

# Appendix

This project utilized multiple data sources, including the SpaceX API for real-time launch data and web scraping to extract additional information from relevant websites. Key tools and libraries included Python for data manipulation, visualization, and model building, SQL for performing exploratory data analysis and querying structured datasets, Folium for creating interactive geospatial maps, and Plotly Dash for building dynamic dashboards to visualize trends and relationships.

Key metrics analyzed throughout the project included launch success rates, average and range of payload capacities, and model performance metrics such as accuracy, precision, recall, and F1-score. These metrics provided valuable insights into the reliability and efficiency of SpaceX's operations. The results and visualizations can be accessed through the project's GitHub repository, which also contains the interactive Folium map showcasing launch sites and their outcomes.

If you need additional details or specific adjustments, let me know!

Thank you!