

Data Science Capstone Project

SPACE X Booter Landing Prediction

Rahul Kashyap



Outline



Executive Summary



Introduction



Methodology



Results



Conclusion



Appendix

EXECUTIVE SUMMARY



Executive Summary



This project focuses on the predictive analysis of SpaceX launch data, aiming to accurately forecast the successful recovery of a rocket's first stage. Through comprehensive data cleaning, exploratory analysis, and interactive visualizations, we've laid the groundwork for our predictive models.



Utilizing Logistic Regression, Support Vector Machine, Decision Tree Classifier, and K Nearest Neighbours, we achieved an accuracy rate of approximately 83.33%. Our findings underscore the significant influence of launch site locations on success rates, providing a strategic foothold for commercial space companies like Space Y. Further refinement and data enrichment are needed to enhance our model's prediction precision.

INTRODUCTION



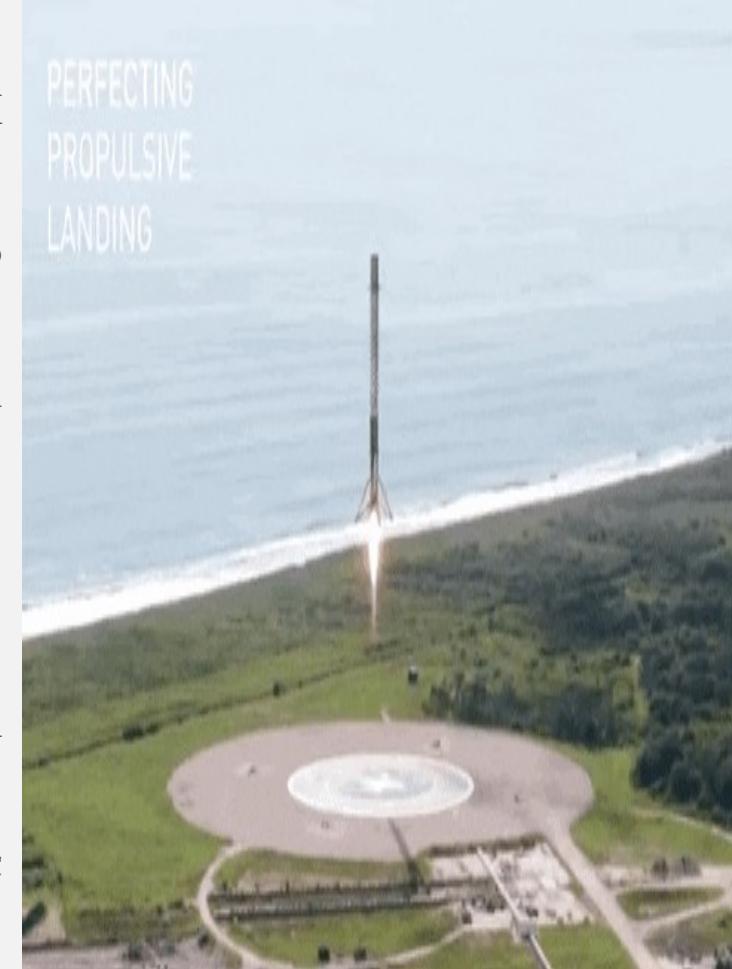
Introduction

Business Understanding:

- The Commercial Space Age demands cost efficiency and sustainability, with Space X leading the charge through successful reuse of the Falcon 9 first stage.
- This strategy allows Space X to undercut competitors drastically (\$62 million vs. \$165 million), making them the industry benchmark.
- Competitors, like Space Y, seek to emulate this success to become competitive and sustainable.

Problem Statement:

- Our task is to develop a machine learning model capable of predicting the successful recovery of the first stage of a launch vehicle of Falcon 9 Rockets
- The final goal is to identify the most effective prediction model, aiding Space Y's strategic planning in the competitive space industry.



METHODOLOGY

Data
Collection

Data
Wrangling

Data
Visualizations

Data ML
Models and
Prediction

DATA COLLECTION



Data Collection

Data collection strategy harnessed the capabilities of both API requests and web scraping to gather comprehensive launch data from public SpaceX resources. The amalgamation of these sources enabled us to curate a robust dataset for subsequent analysis and modelling.

- The SpaceX API offered us a multitude of attributes including FlightNumber, BoosterVersion, PayloadMass, LaunchSite, Outcome, and more, while web scraping SpaceX's Wikipedia page supplemented our dataset with additional parameters such as Launch site, Customer, Launch outcome, and Booster landing.
- Ensuring data integrity and handling missing data were critical challenges during this phase, but they were meticulously addressed, considering the pivotal role accurate data collection plays in ensuring the reliability of our analysis and predictive models.

# Hint data['BoosterVersion']!='Falcon 1'													
data_falcon9 = data2[data2['BoosterVersion']=='Falcon 9']													
data_falcon9													
FlightNumber	Date	BoosterVersion	PayloadMass	Orbit	LaunchSite	Outcome	Flights	GridFins	Reused	Legs	LandingPad	Block	R
4	6 2010-06-04	Falcon 9	Nan	LEO	CCSFS SLC 40	None None	1 False	False	False	None	None	1.0	
5	8 2012-05-22	Falcon 9	525.0	LEO	CCSFS SLC 40	None None	1 False	False	False	None	None	1.0	
6	10 2013-03-01	Falcon 9	677.0	ISS	CCSFS SLC 40	None None	1 False	False	False	None	None	1.0	
7	11 2013-09-29	Falcon 9	500.0	PO	VAFB SLC 4E	False Ocean	1 False	False	False	None	None	1.0	
8	12 2013-12-03	Falcon 9	3170.0	GTO	CCSFS SLC 40	None None	1 False	False	False	None	None	1.0	

2020 [edit]									
In late 2019, Gwynne Shotwell stated that SpaceX hoped for as many as 24 launches for Starlink satellites in 2020. ^[490] In addition to 14 or 15 non-Starlink launches, At 26 launches, 13 of which for Starlink satellites, Falcon 9 had its most prolific year, and Falcon rockets were second most prolific rocket family of 2020, only behind China's Long March rocket family. ^[491]									
Flight No.	Date and time (UTC)	Version, Booster ^[b]	Launch site	Payload ^[d]	Payload mass	Orbit	Customer	Launch outcome	Booster landing
78	02:19:21 ^[492] 7 January 2020,	F9 B5 Δ B1049.4	CCAFS, SLC-40	Starlink 2 v1.0 (60 satellites)	15,600 kg (34,400 lb) ^[8]	LEO	SpaceX	Success	Success (drone ship)
	14:07 ^[493] 19 January 2020, 15:30 ^[494]	F9 B5 Δ B1046.4	KSC, LC-39A	Crew Dragon in-flight abort test ^[495] (Dragon C205.1)	12,050 kg (26,570 lb)	Sub-orbital ^[496]	NASA (CTS) ^[497]	Success	No attempt
79	14:07 ^[498] 29 January 2020, 15:05 ^[500]	F9 B5 Δ B1051.3	CCAFS, SLC-40	Starlink 2 v1.0 (60 satellites)	15,600 kg (34,400 lb) ^[8]	LEO	SpaceX	Success	Success (drone ship)
	17 February 2020, 04:59 ^[501]	F9 B5 Δ B1056.4	CCAFS, SLC-40	Starlink 4 v1.0 (60 satellites)	15,600 kg (34,400 lb) ^[8]	LEO	SpaceX	Success	Failure (drone ship)
80	14:07 ^[502] 04:59 ^[503] 7 March 2020, 04:59 ^[504]	F9 B5 Δ B1059.2	CCAFS, SLC-40	SpaceX CRS-20 (Dragon C112.3 Δ)	1,977 kg (4,359 lb) ^[505]	LEO (ISS)	NASA (CRS)	Success	Success (ground pad)
	18 March 2020, 12:16 ^[510]	F9 B5 Δ B1048.5	KSC, LC-39A	Starlink 5 v1.0 (60 satellites)	15,600 kg (34,400 lb) ^[8]	LEO	SpaceX	Success	Failure (drone ship)
81	14:07 ^[511] 22 April 2020, 19:30 ^[514]	F9 B5 Δ B1051.4	CCAFS, SLC-40	Starlink 6 v1.0 (60 satellites)	15,600 kg (34,400 lb) ^[8]	LEO	SpaceX	Success	Success (drone ship)

[GitHub Link 1](#)

[GitHub Link 2](#)

DATA WRANGLING



Data Wrangling

Data Wrangling process was conducted with a focus on ensuring the cleanliness, consistency, and utility of our data set. This involved handling missing values, removing duplicates, and converting data types, with the ultimate aim of preparing our data for accurate analysis and predictive modelling.

Key steps and outcomes:

- **Training Label Creation:** We constructed a new training label column, 'class', based on the 'Mission Outcome' and 'Landing Location'. Successful landings (True ASDS, True RTLS, True Ocean) were encoded as '1', while unsuccessful or unknown outcomes (None None, False ASDS, None ASDS, False Ocean, False RTLS) were encoded as '0'. This binary classification sets the stage for our machine learning models.
- **Data Preprocessing:** Utilizing Python's pandas library, we meticulously handled missing values, removed duplicates, and converted data types to appropriate formats. This yielded a well-structured dataset ready for exploratory analysis and modelling.
- **Challenges and Resolution:** Challenges during data wrangling included dealing with missing or inconsistent data and deciding on how to handle outliers. These were systematically addressed, ensuring the final dataset is reliable and robust for downstream applications.

The importance of data wrangling in our process cannot be overstated, as it is a crucial step that significantly affects the accuracy of our subsequent analysis and models.

[GitHub Link](#)

DATA VISUALIZATION Methodology



Methodology - Exploratory Data Analysis (EDA)

Exploratory Data Analysis (EDA) was instrumental in uncovering the underlying structure of our data, identifying outliers, checking assumptions, and revealing vital patterns and relationships. This comprehensive analysis informed our subsequent modeling process, ensuring our insights were data-driven and robust.

Key aspects and findings:

- **Variables Analyzed:** The primary focus of our EDA was on variables such as Flight Number, Payload Mass, Launch Site, Orbit, Class, and Year. These were identified as significant potential influencers on the successful recovery of a rocket's first stage.
- **Visualization Tools:** Utilizing Python's seaborn library, we created a range of visual representations like scatter plots, line charts, and bar plots to help unravel the relationships between the selected variables.
- **Key Visualizations:** We examined the relationships via plots of Flight Number vs. Payload Mass, Flight Number vs. Launch Site, Payload Mass vs. Launch Site, Orbit vs. Success Rate, Flight Number vs. Orbit, Payload vs Orbit, and Success Yearly Trend. These visualizations facilitated our understanding of which variables could significantly influence the outcomes and should therefore be used in our machine learning model training.
- **Challenges:** Our EDA encountered challenges such as dealing with large datasets and visualizing high-dimensional data. These were tackled by using appropriate plotting methods and focusing on key influential variables.

The insights derived from the EDA were critical in shaping our modeling strategy and enhancing the predictive capability of our machine learning models.

Methodology - Exploratory Data Analysis (EDA) (SQL)

The Exploratory Data Analysis (EDA) phase, underpinned by **SQL**, was aimed at revealing the underlying structure of our data, validating assumptions, and pinpointing outliers and potential anomalies. It offered us an understanding of the data that informed our subsequent modeling strategy and highlighted key patterns and relationships.

Key aspects and findings:

- **SQL Integration:** We leveraged SQL's querying power through Python to delve into our dataset, which was loaded into an IBM DB2 Database. This strategy allowed for seamless interaction and in-depth examination of the data.
- **Variables and Relationships Explored:** SQL queries enabled us to extract detailed information about launch site names, mission outcomes, various payload sizes of customers, booster versions, and landing outcomes. The analysis illuminated essential relationships between variables, particularly between launch site and success rate.
- **Visualization Tools:** Complementing SQL, we used Python's seaborn library for visualizing our data, translating the queried data into insightful plots. This combination provided a robust exploratory framework, enabling us to tackle challenges associated with large and high-dimensional datasets.
- **Outcome and Influence:** The insights gleaned from the SQL-based EDA significantly influenced our modeling process. It informed our understanding of potential feature importance and guided our feature selection for the machine learning model training.

The EDA, enriched with SQL, served as a significant step in our data analysis pipeline, enhancing the rigor and comprehensiveness of our approach.

[GitHub Link](#)

Methodology - Interactive Visual Analytics using Folium and Plotly Dash

Our project extensively leverages the power of interactive visual analytics, utilizing tools like Folium and Plotly Dash to create engaging visualizations. This enables detailed data exploration and brings forth deeper insights, helping us understand the nuances of launch site locations and outcomes.

Key aspects and findings:

- **Use of Folium:** We used Folium to create interactive maps marking launch sites, successful and unsuccessful landings, and their proximity to key locations such as railways, highways, coasts, and cities. This visual aid enhanced our understanding of the strategic placement of launch sites and the relationship between location and successful recoveries.
- **Use of Plotly Dash:** Our dashboard, created using Plotly Dash, included a pie chart and a scatter plot. The pie chart could be tailored to show the distribution of successful landings across all launch sites or individual launch site success rates. The scatter plot, taking launch site and payload mass as inputs, showed variations in success rates across different parameters.
- **Importance of Interactive Visual Analytics:** These visualizations were crucial in allowing a detailed, engaging exploration of the data. They enhanced our understanding of factors influencing successful recoveries and informed our subsequent modeling process.
- **Challenges:** The process of creating user-friendly, interactive visualizations with large datasets posed certain challenges. However, careful selection of visualization tools and appropriate data handling techniques helped overcome these hurdles.

Interactive visual analytics significantly enriched our data exploration process, providing insightful visuals that informed our modeling and prediction approach.

[GitHub Link](#)

DATA MODEL



Data Model - Classification & Prediction

- **Overview of Predictive Analysis:** Using classification models to predict launch outcomes. The goal is to use the historical data to predict future outcomes, specifically the success of SpaceX launches.
- **Data Split:** The data was split into a training set and a test set with a **test size of 20%**. The training set is used to train the models, and the test set is used to evaluate their performance. This helps to ensure that the models are able to generalize well to new, unseen data.
- **Standardization:** Before training the models, the data was standardized using the StandardScaler from sklearn. This step is crucial because many machine learning algorithms do not perform well if the features are not on the same scale. StandardScaler standardizes the features by removing the mean and scaling to unit variance.
- **Models Used:** Several classification models were used in the analysis, including Support Vector Machines (SVM), Classification Trees, Logistic Regression, and K-Nearest Neighbors (KNN). Each of these models has its own strengths and weaknesses, and they are suited to different types of data and tasks.
 - **SVM:** SVM is a powerful and flexible classification algorithm that can handle both linear and non-linear data. It works by finding the hyperplane that best separates the classes in the data.
 - **Classification Trees:** Classification trees are simple yet effective models that make predictions based on a series of decision rules. They are easy to understand and interpret, making them a good choice for exploratory analysis.
 - **Logistic Regression:** Logistic regression is a simple and fast model that is particularly well-suited to binary classification tasks. It works by fitting a logistic function to the data.
 - **KNN:** KNN is a non-parametric method that makes predictions for a new observation by searching the entire dataset for the K observations that are closest to it.
- **Hyperparameter Tuning:** For each model, the best hyperparameters were found using GridSearchCV. This is a process of performing hyperparameter tuning in order to determine the optimal values for a given model. This is significant as the performance of the entire model is based on the hyperparameter values specified.
- **Model Evaluation:** The models were evaluated using the accuracy score, which is the proportion of correct predictions out of the total predictions. The accuracy scores were compared to determine the best performing model.

RESULT



Result – Data Wrangling

```
df['LaunchSite'].value_counts()  
  
CCAFS SLC 40      55  
KSC LC 39A        22  
VAFB SLC 4E       13  
Name: LaunchSite, dtype: int64
```

```
df['Reused'].value_counts()  
  
False    53  
True     37  
Name: Reused, dtype: int64
```

```
# landing_outcomes = values on Outcome column  
landing_outcomes = df['Outcome'].value_counts()  
landing_outcomes
```

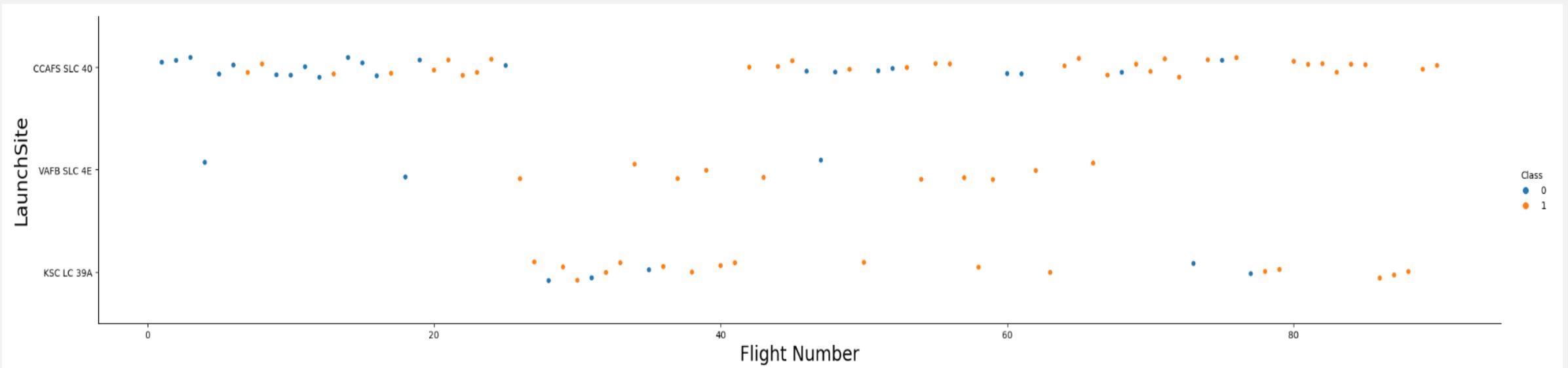
```
True ASDS      41  
None None     19  
True RTLS     14  
False ASDS    6  
True Ocean    5  
False Ocean   2  
None ASDS     2  
False RTLS    1  
Name: Outcome, dtype: int64
```

```
# Apply value_counts on Orbit column  
df['Orbit'].value_counts()
```

```
GT0      27  
ISS      21  
VLEO     14  
PO       9  
LEO      7  
SSO      5  
MEO      3  
ES-L1    1  
HEO      1  
SO       1  
GEO      1  
Name: Orbit, dtype: int64
```

- ✓ There are 3 different launch sites with most launch being from CCAFS SLC 40 (i.e. Cape Canaveral Space Launch Complex)
- ✓ Reused percentage being at 41% as of now
- ✓ True ASDS means the mission outcome was successfully landed to a drone ship
- ✓ GTO is geosynchronous orbit is a high Earth orbit that allows satellites to match Earth's rotation. Located at 22,236 miles (35,786 kilometers) above Earth's equator, this position is a valuable spot for monitoring weather, communications and surveillance

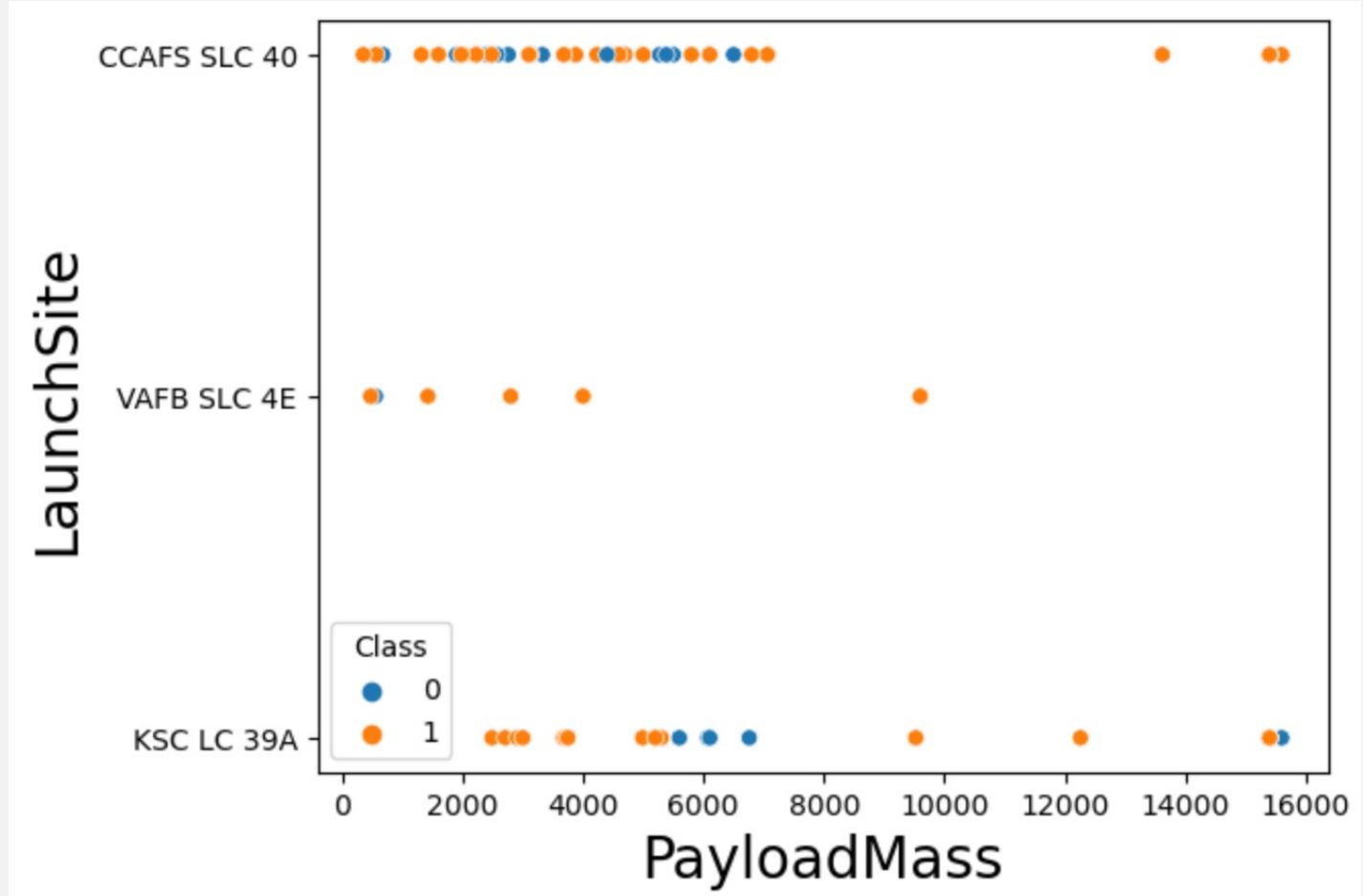
Result – EDA (Launch Site vs Flight Number)



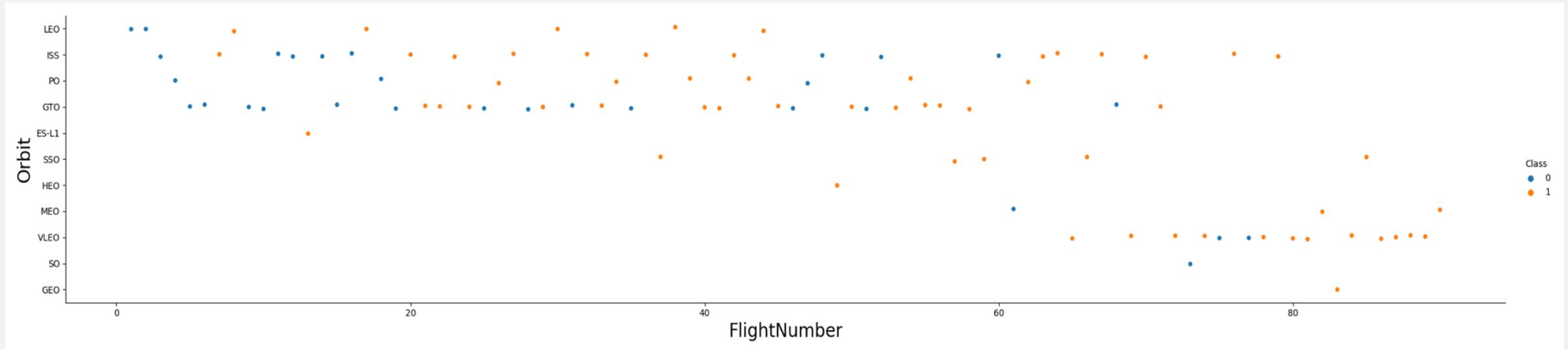
- ✓ Over time, we have observed an increased success rate in launches, with a significant breakthrough happening around the 20th flight. This increase in success rate indicates advancements in technology and experience gained over the iterations.
- ✓ The most active and successful launch site is Cape Canaveral Air Force Station (CCAFS SLC 40), accounting for the highest volume of successful launches.
- ✓ Following CCAFS, the next successful launch sites are Vandenberg Air Force Base (VAFB SLC 4E) and Kennedy Space Center (KSC LC 39A).
- ✓ The overall success rate across all sites has improved over time, reflecting continuous learning and improvements in launch processes.

Result - EDA (Launch Site vs Payload)

- ✓ The majority of payload masses range between 0-6000 kg across all launch sites.
- ✓ Payloads exceeding 9,000 kg, approximately the weight of a school bus, demonstrate a notably high success rate.
- ✓ Heavier payloads over 12,000 kg are predominantly launched from two sites: CCAFS SLC 40 and KSC LC 39A.
- ✓ The ability to successfully launch heavier payloads could be a unique capability of these two sites, suggesting a potential specialization in handling high-mass launches.

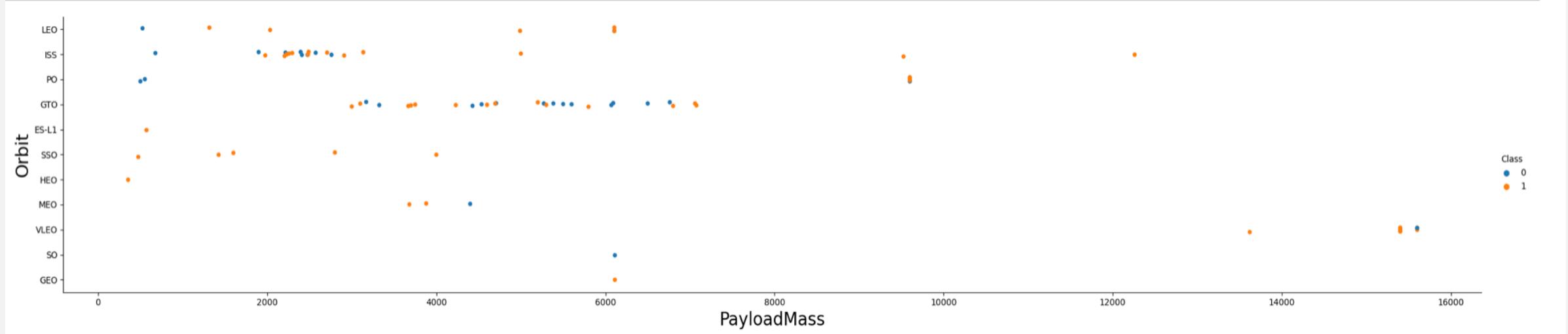


Result – EDA (Orbit Type vs Flight Number)



- ✓ SpaceX's launch orbit preferences evolved over time, with a noticeable shift in success rates correlating with these changes.
- ✓ Initially, SpaceX primarily targeted Low Earth Orbit (LEO) which saw moderate success rates.
- ✓ Over time, SpaceX expanded its focus to include Very Low Earth Orbit (VLEO), which has been associated with an increase in successful launches.
- ✓ The success rate across all orbits has shown a general improvement over time, indicating advancements in SpaceX's launch capabilities.
- ✓ The recent increase in VLEO launches suggests that SpaceX is exploring new business opportunities in this orbit.

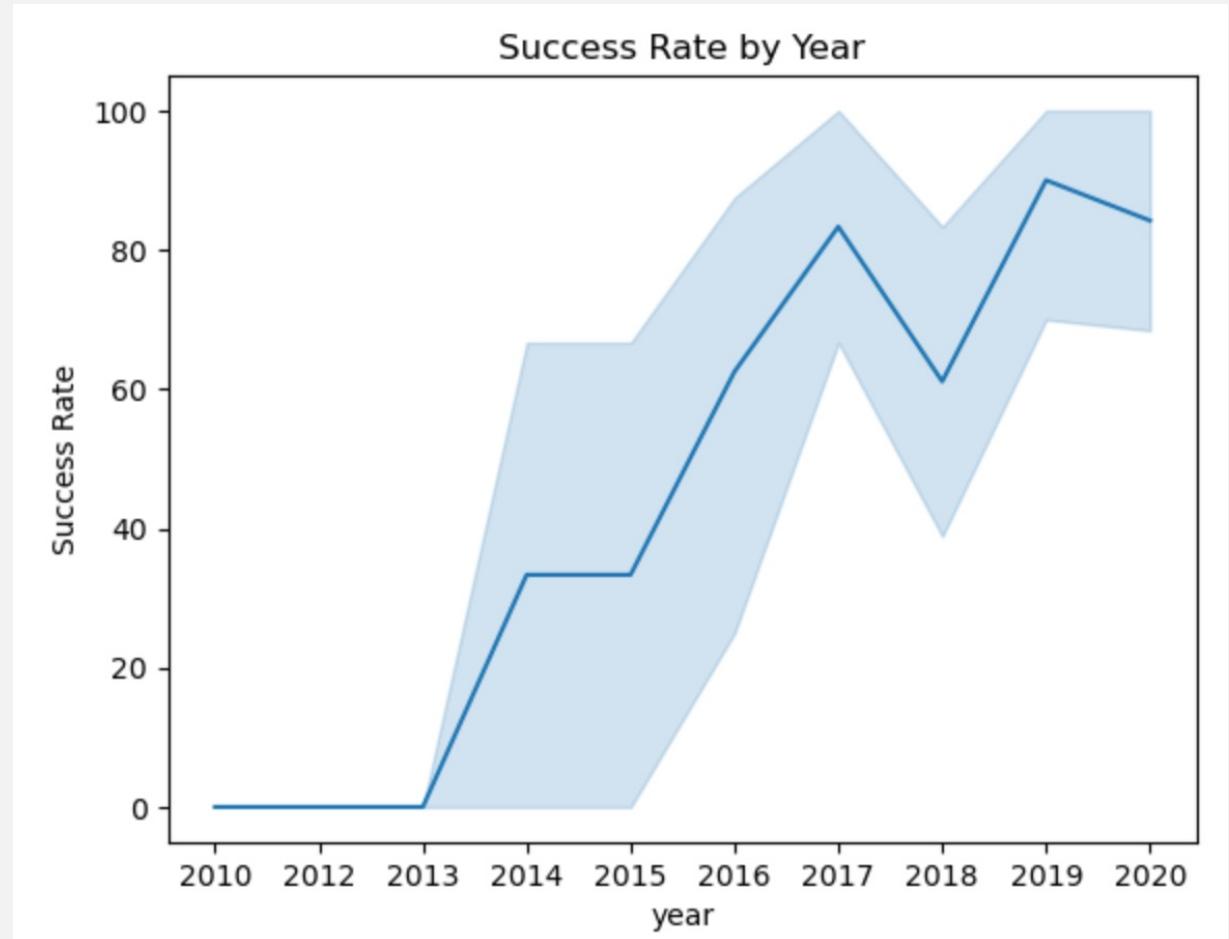
Result – EDA (Orbit Type vs Payload Mass)



- ✓ The analysis reveals no apparent correlation between payload mass and success rate for Geosynchronous Transfer Orbit (GTO).
- ✓ The International Space Station (ISS) orbit displays a wide range of payload masses, with a generally high success rate.
- ✓ Sun-Synchronous Orbit (SSO) and Low Earth Orbit (LEO) typically accommodate lower payload masses.
- ✓ Very Low Earth Orbit (VLEO), another orbit with high success rates, predominantly handles higher payload masses.
- ✓ Few launches have been directed towards Semi-Synchronous Orbit (SO) and Geosynchronous Orbit (GEO), limiting the analysis for these orbits.
- ✓ Overall, the success of a launch appears to be influenced by a combination of the orbit type and the payload mass.

Result - EDA (Launch Success Yearly Trend)

- ✓ The success rate of SpaceX launches has shown a general upward trend since 2013, indicating continuous improvements in technology and operational processes.
- ✓ The initial years (2010-2012) were marked by a lower success rate, likely due to the company's nascent stage and the challenges of developing new rocket technology.
- ✓ A significant increase in success rate was observed starting from 2013, marking a period of technological maturity and consistent performance.
- ✓ Despite the overall positive trend, a slight dip in success rate was observed in 2018, which could be attributed to specific challenges or anomalies encountered that year.
- ✓ In recent years, the success rate has stabilized at around 80%, showcasing SpaceX's consistent performance in successful launches.
- ✓ The 95% confidence interval indicates a high level of certainty in these trends, with the light blue shading representing the range within which the true success rate lies with 95% confidence.

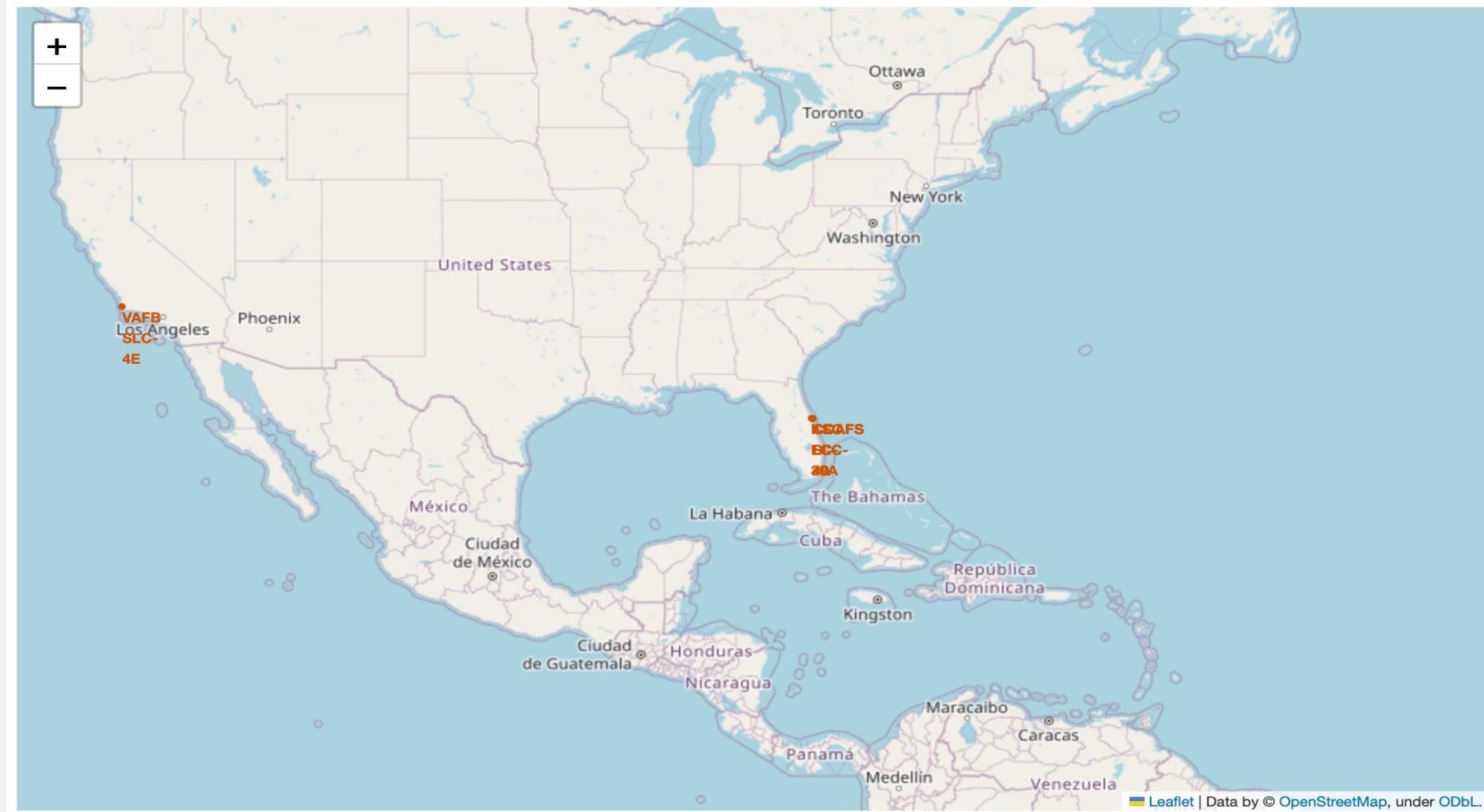


Result – Ranking Counts of Successful Landings

Ranking of all landing outcomes between the date 2010-06-04 and 2017-03-20

Landing Outcome	Occurrences
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

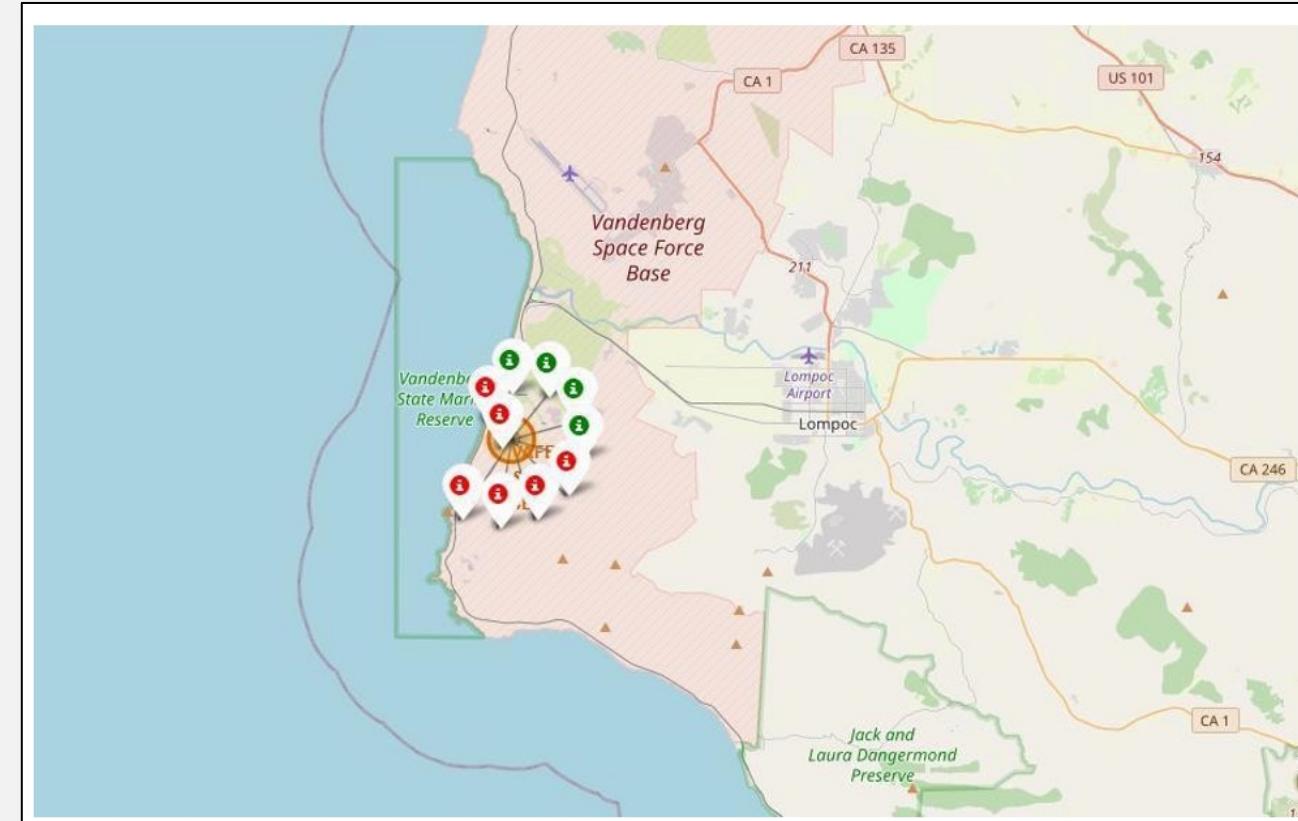
Result – Folium Map Visualization(Launch Site)



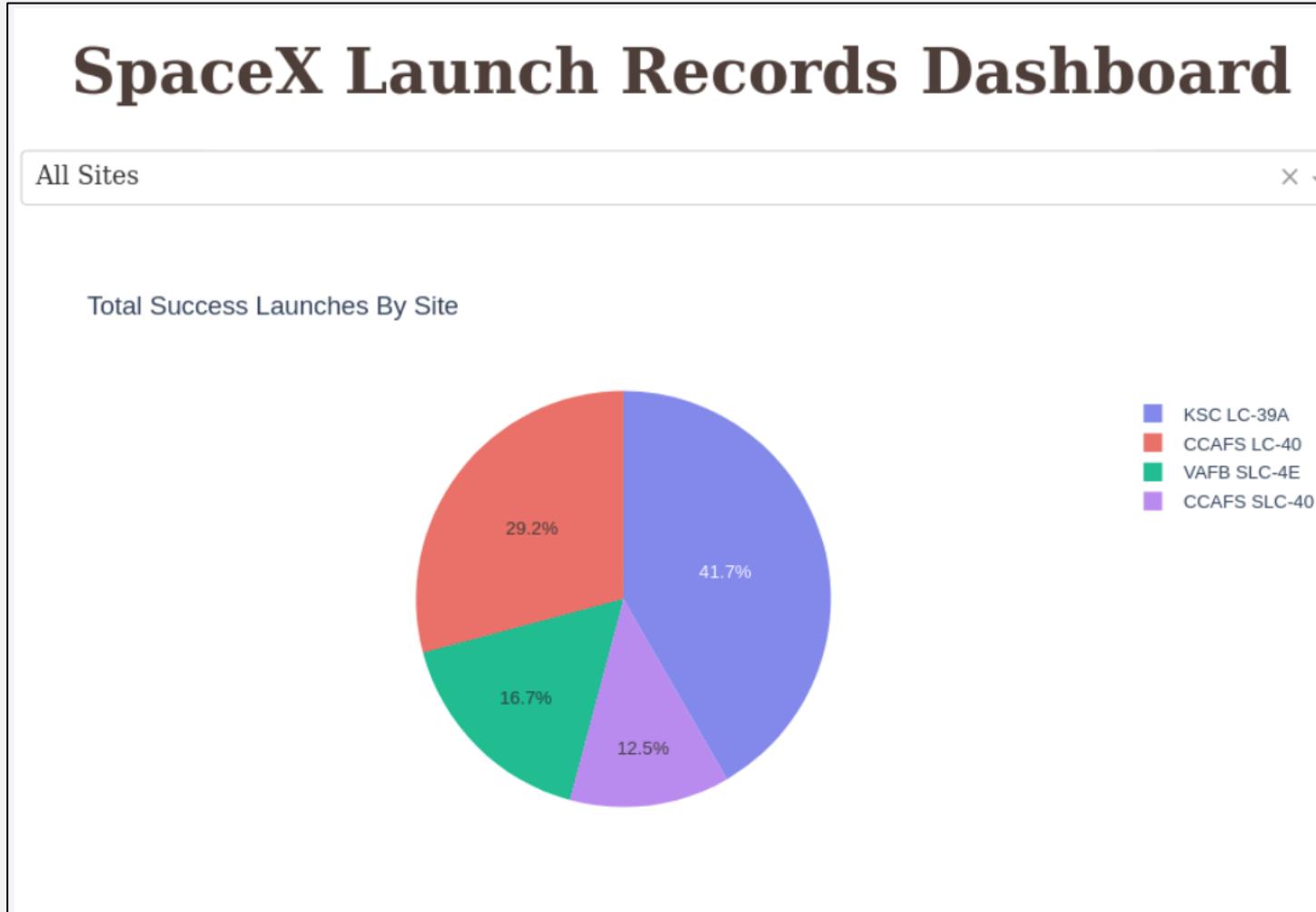
Result - Folium Map Visualization(Launch Outcome)

Example of VAFB SLC-4E launch site launch outcomes

Green markers indicate successful and red ones indicate failure



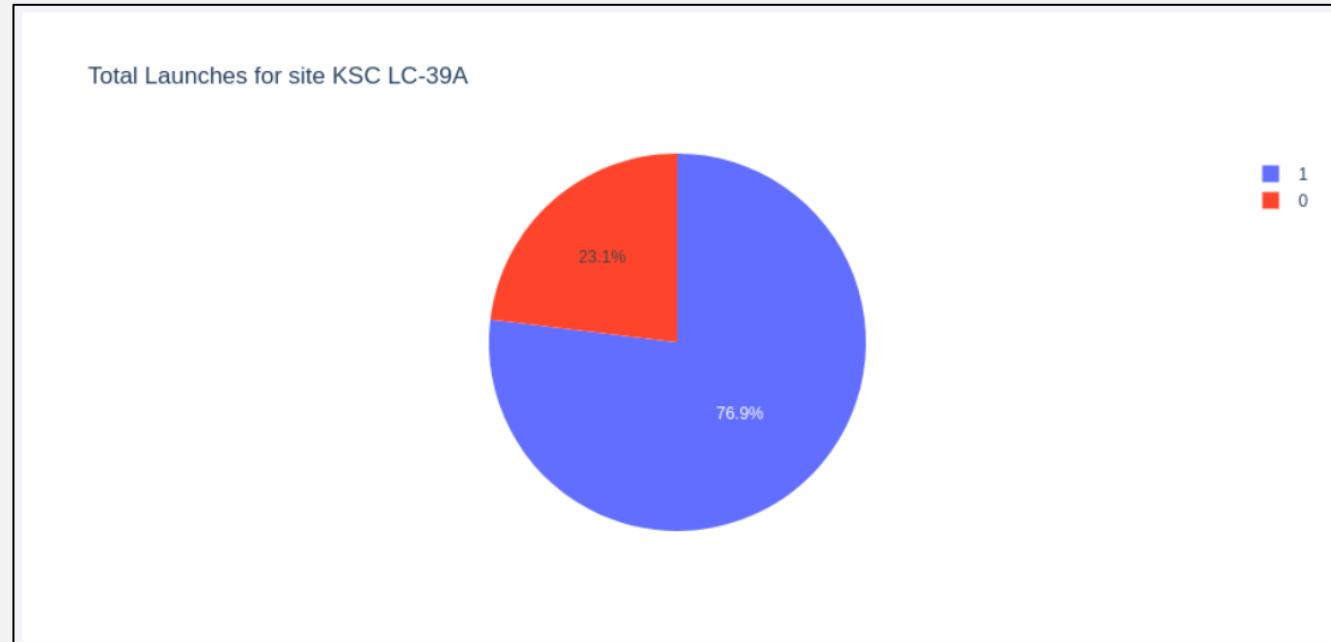
Result – Plotly Dashboard



- ✓ The distribution of successful landings across all launch sites shows that CCAFS LC-40 (old name) and CCAFS SLC-40 (new name) have the highest number of successful landings.
- ✓ The name change does not affect the success rate, indicating consistent performance at this site.
- ✓ VAFB SLC-4E has the lowest share of successful landings, potentially due to a smaller number of launches and increased difficulty of launching from the west coast.
- ✓ The launch site appears to be a significant factor in the success of missions, suggesting the importance of geographical and logistical considerations in planning launches.

Result – Plotly Dashboard

- ✓ The Plotly Dashboard analysis reveals that KSC LC-39A is the launch site with the highest success rate.
- ✓ This site has a total of 13 launches, out of which 10 were successful and 3 were unsuccessful.
- ✓ The success rate at KSC LC-39A is therefore approximately 76.9%, indicating a high level of reliability for launches from this site.
- ✓ This high success rate could be attributed to various factors such as location advantages, technical expertise, and advanced infrastructure, which could be explored in further analyses.
- ✓ The findings suggest that focusing on successful launch sites like KSC LC-39A could potentially improve overall mission success rates for SpaceX.

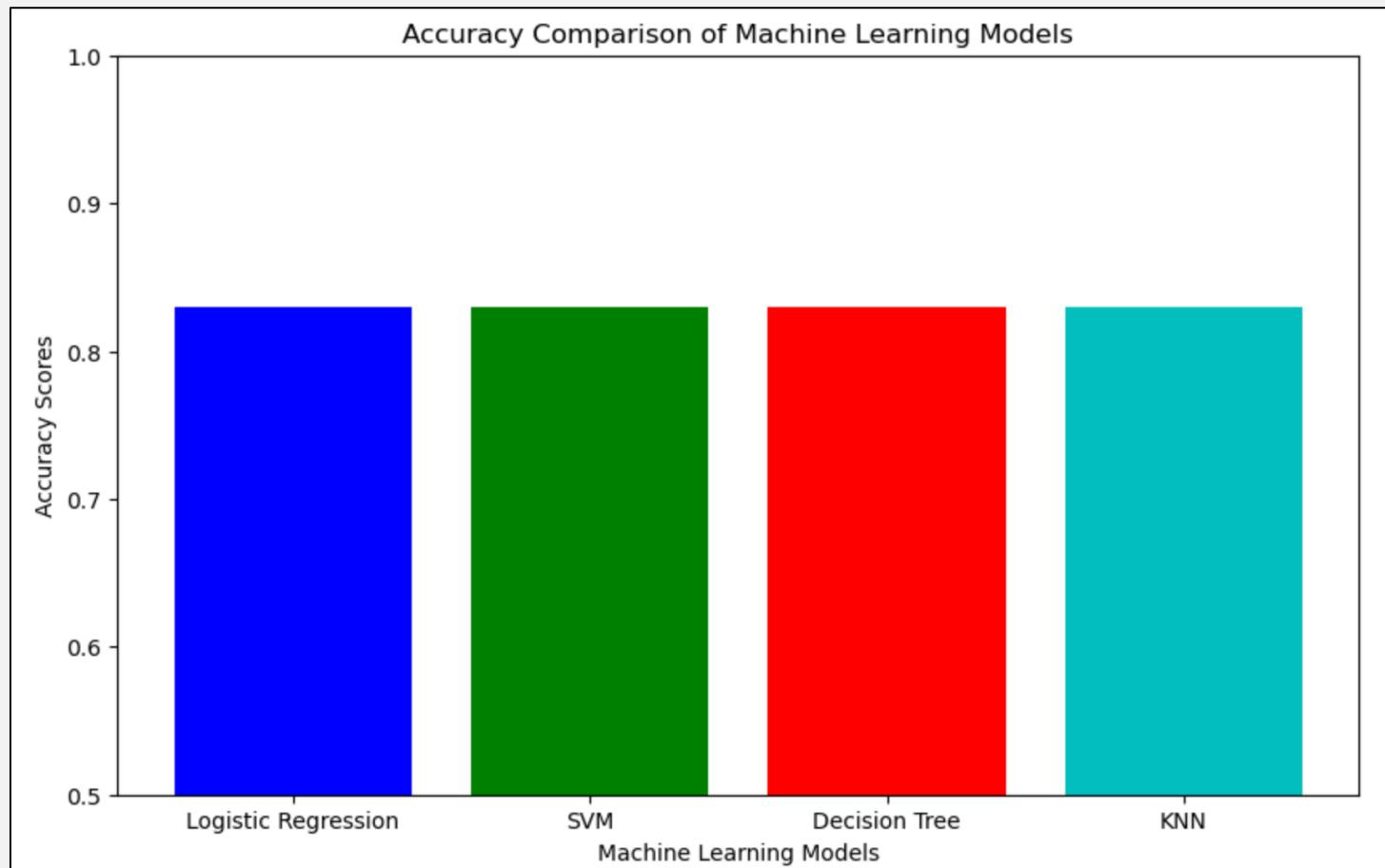


RESULT - DATA MODELING AND PREDICTION



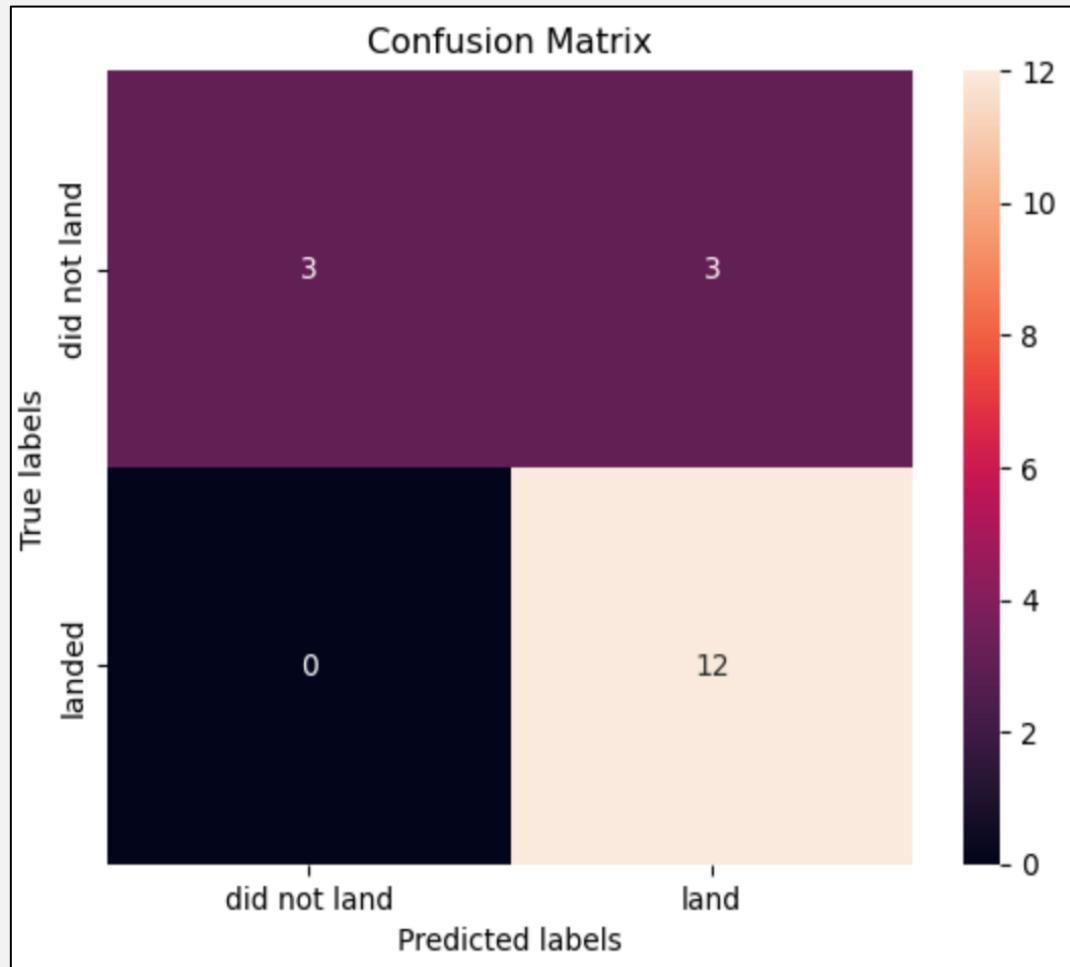
Result – Classification Model Accuracy

- ✓ Four classification models were used: Logistic Regression, Support Vector Machine (SVM), Decision Tree, and K-Nearest Neighbors (KNN).
- ✓ All models, except for the Decision Tree Classifier, achieved a similar accuracy of 83.33% on the test set.
- ✓ The Decision Tree Classifier outperformed other models with an accuracy of over 87%.
- ✓ However, it's important to note that the test size was small (only 18 samples), which can cause large variance in accuracy results, especially noticeable in the Decision Tree Classifier model in repeated runs.
- ✓ Due to the small sample size, more data might be needed to definitively determine the best model for predicting launch success.



Result – Confusion Matrix

- ✓ All models (Logistic Regression, SVM, KNN) predicted 12 successful landings correctly when the actual outcome was a successful landing
- ✓ The models correctly predicted 3 unsuccessful landings when the actual outcome was an unsuccessful landing.
- ✓ However, the models incorrectly predicted 3 successful landings when the actual outcome was an unsuccessful landing, indicating a tendency to over-predict successful landings (false positives)
- ✓ The Decision Tree Classifier demonstrated its accuracy through the confusion matrix, with a higher number of true positives and true negatives compared to false positives and false negatives
- ✓ Despite the over-prediction of successful landings, the models still demonstrated a high level of accuracy in their predictions
- ✓ The confusion matrix results highlight the importance of further refining the models to reduce the number of false positives, potentially improving the prediction accuracy for unsuccessful landings



CONCLUSION



Conclusion



The project successfully developed a machine learning model with an accuracy of 83% to predict the success of SpaceX's Stage 1 landing, potentially saving around \$100 million per launch.



The best launch site for successful landings was identified as KSC LC-39A, and launches with a payload above 7,000kg were found to be less risky.



The analysis revealed an improvement in successful landing outcomes over time, indicating the evolution and enhancement of SpaceX's processes and rockets.



The use of interactive visual analytics tools, Folium and Plotly Dash, enabled a detailed and engaging exploration of the data, uncovering key patterns and relationships.



While the machine learning models demonstrated high accuracy, there is room for further refinement and improvement, particularly in reducing the number of false positives.



The findings from this analysis can significantly inform SpaceX's launch strategy, potentially contributing to improved launch success rates and increased profits in the future.



For future work, it is recommended to collect more data to further refine the machine learning models and improve prediction accuracy.

Innovative Insights



The use of Folium for creating interactive maps of launch sites and launch outcomes is an innovative approach to visualizing geographical data.



The use of Plotly Dash for building a dashboard for detailed launch records provides an interactive and user-friendly way to explore the data.



The application of various machine learning models for predicting launch outcomes is a novel use of predictive analytics in the context of space launches.



The analysis revealed that launch site location has a significant influence on launch success rate, an insight that could have far-reaching implications for SpaceX's launch strategy.

Appendix

- Github link: https://github.com/Rahul-Kashyap2/DataScience_Capstone_Project

THANK YOU !

