

Winning Space Race with Data Science

Akhil Nishad
12-June-2025



Outline

- Executive Summary
- Introduction
- Methodology
- Results
- Conclusion

Executive Summary

- Summary of methodologies
 - Data Collection through API
 - Data Collection with Web Scraping
 - Data Wrangling
 - Exploratory Data Analysis with SQL
 - Exploratory Data Analysis with Data Visualization
 - Interactive Visual Analytics with Folium
 - Machine Learning Prediction
- Summary of all results
 - Exploratory Data Analysis result
 - Interactive analytics in screenshots
 - Predictive Analytics result

Introduction

- **Project background and context**

SpaceX reduces launch costs by reusing the Falcon 9 first stage. Predicting whether the first stage will land helps estimate launch costs. In this project, as a data scientist at fictional company SpaceY, I use public SpaceX data and machine learning to predict landings and support business decisions through dashboards.

- **Problems you want to find answers**

- Can we predict if the first stage will land?
- What factors affect landing success?
- How does launch performance vary over time?
- How can we visualize insights for decision-making?

Section

1

Methodology

Methodology

Executive Summary

- Data collection methodology:
 - Data was collected using the SpaceX API and web scraping from wikipedia
- Perform data wrangling
 - One hot encoding was applied to categorical features
- 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

- The data was collected using various methods
 - Data collection was done using get request to the SpaceX API.
 - Next, we decoded the response content as a Json using `.json()` function call and turn it into a pandas dataframe using `.json_normalize()`.
 - We then cleaned the data, checked for missing values and fill in missing values where necessary.
 - In addition, we performed web scraping from Wikipedia for Falcon 9 launch records with BeautifulSoup.
 - The objective was to extract the launch records as HTML table, parse the table and convert it to a pandas dataframe for future analysis.

Data Collection – SpaceX API

- We retrieved data from the SpaceX API using a GET request, then cleaned, formatted, and performed basic wrangling on the data to prepare it for analysis.

- Github url: [link](#)

The screenshot shows a Jupyter Notebook interface with several code cells and explanatory text.

```
spaces_url="https://api.spacexdata.com/v4/launches/past"
[1]: ✓ 0.0s Python
response = requests.get(spaces_url)
[1]: ✓ 1.6s
Check the content of the response
[2]: ✓ 0.0s Python
print(response.content)
[2]: ✓ 0.0s
b'[{\"fairings\":{\"reused\":false,\"recovery_attempt\":false,\"recovered\":false,\"ships\":[]},\"links\":{\"patch\":{\"small\":\"https://images2.imgur.com/94/f2/WNGPh45r_o.png\",\"large\":\"https://images2.imgur.com/5b/82/QcxHbBd...}'}
You should see the response contains massive information about SpaceX launches. Next, let's try to discover some more relevant information for this project.

Task 1: Request and parse the SpaceX launch data using the GET request
To make the requested JSON results more consistent, we will use the following static response object for this project:
```

```
static_json_url='https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBN-058321BN-SkillNetwork/datasets/API_call_spacex_api.json'
[3]: ✓ 0.0s Python
response = requests.get(static_json_url)
[3]: ✓ 0.0s
We should see that the request was successful with the 200 status response code
[4]: ✓ 0.0s Python
response.status_code
[4]: ✓ 0.0s
200
Now we decode the response content as a Json using .json() and turn it into a Pandas dataframe using .json_normalize()
[5]: ✓ 0.0s Python
# Use json_normalize method to convert the json result into a dataframe
data = pd.json_normalize(response.json())
[5]: ✓ 0.0s
```

Data Collection - Scraping

- We used web scraping with BeautifulSoup to extract Falcon 9 launch records from the web. The data was parsed from an HTML table and converted into a pandas DataFrame for analysis.
- Github url: [Link](#)

The screenshot shows a Jupyter Notebook interface with several code cells and explanatory text.

Code Cell 1:

```
static_url = "https://en.wikipedia.org/w/index.php?title=List_of_Falcon_9_and_Falcon_Heavy_launches&oldid=1827686922"
```

Next, request the HTML page from the above URL and get a `response` object

TASK 1: Request the Falcon9 Launch Wiki page from its URL

First, let's perform an HTTP GET method to request the Falcon9 Launch HTML page, as an HTTP response.

```
# use requests.get() method with the provided static_url  
# assign the response to a object  
response = requests.get(static_url)
```

Code Cell 2:

```
# use BeautifulSoup() to create a BeautifulSoup object from a response text content  
soup = BeautifulSoup(response.text)
```

Create a BeautifulSoup object from the HTML response

Code Cell 3:

```
# Use soup.title attribute  
soup.title
```

Print the page title to verify if the BeautifulSoup object was created properly

```
... <title>List of Falcon 9 and Falcon Heavy launches – Wikipedia</title>
```

TASK 2: Extract all column/variable names from the HTML table header

Next, we want to collect all relevant column names from the HTML table header

Let's try to find all tables on the wiki page first. If you need to refresh your memory about BeautifulSoup, please check the external reference link towards the end of this lab

```
html_tables = soup.find_all('table')
```

Code Cell 4:

```
[1] ✓ 0.0s
```

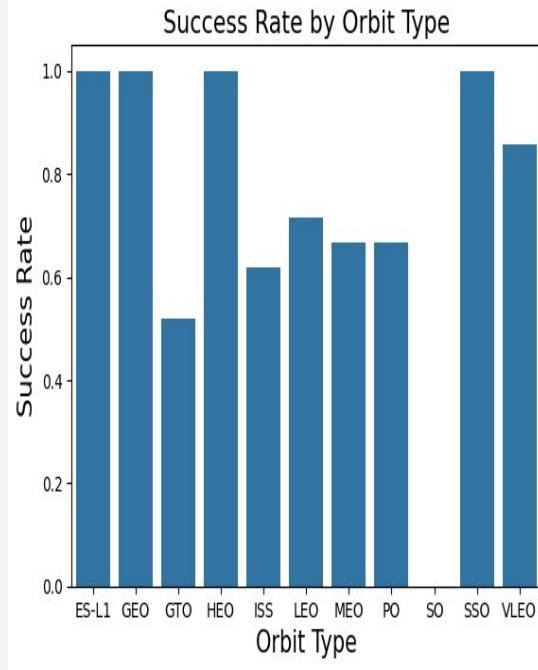
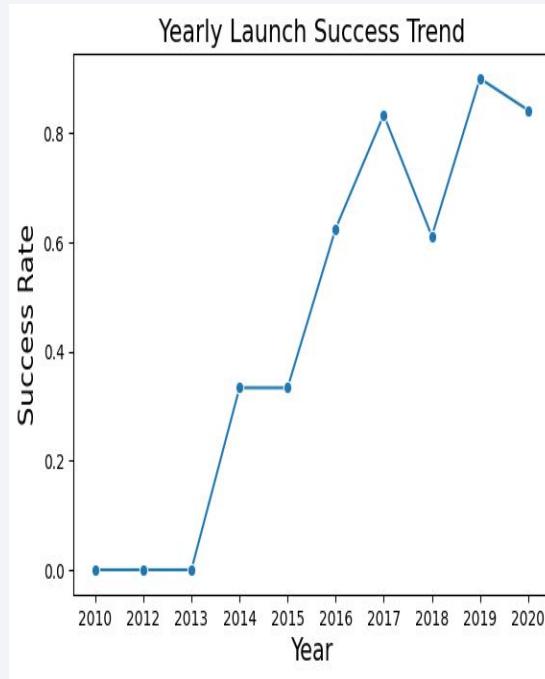
Python

Data Wrangling

- We conducted exploratory data analysis and defined the training labels.
- This included analyzing the number of launches per site and the frequency of different orbit types.
- We also derived the landing outcome label from the outcome column and exported the processed data to a CSV file.
- Github url: [Link](#)

EDA with Data Visualization

- We visualized the data to explore relationships such as flight number vs. launch site, payload vs. launch site, success rates by orbit type, flight number vs. orbit type, and launch success trends over the years.
- Github url: [Link](#)



EDA with SQL

- We loaded the SpaceX dataset into a SQLite database directly within the Jupyter notebook. Using SQL, we performed exploratory analysis to uncover insights such as:
 - Identifying unique launch site names
 - Calculating total payload mass for NASA (CRS) missions
 - Finding the average payload mass for booster version F9 v1.1
 - Counting total successful and failed mission outcomes
 - Retrieving details of failed drone ship landings, including booster versions and launch sites.
- Github url: [Link](#)

Build an Interactive Map with Folium

- We visualized all launch sites on a Folium map, adding markers, circles, and lines to indicate launch outcomes.
- Launch outcomes were classified as 0 for failure and 1 for success.
- By using color-coded marker clusters, we identified launch sites with higher success rates.
- Additionally, we calculated distances from each launch site to nearby locations.
- Github url: [Link](#)

Build a Dashboard with Plotly Dash

- We developed an interactive dashboard using Plotly Dash.
- It includes pie charts displaying the total number of launches by site, and scatter plots illustrating the relationship between launch outcomes and payload mass (kg) across different booster versions.
- Github url: [Link](#)

Predictive Analysis (Classification)

- We used NumPy and pandas to load and preprocess the data, then split it into training and testing sets.
- Multiple machine learning models were built and optimized using GridSearchCV for hyperparameter tuning.
- Accuracy was used as the evaluation metric, and model performance was enhanced through feature engineering and algorithm optimization.
- Finally, we identified the best-performing classification model.
- Github url: [Link](#)

Results

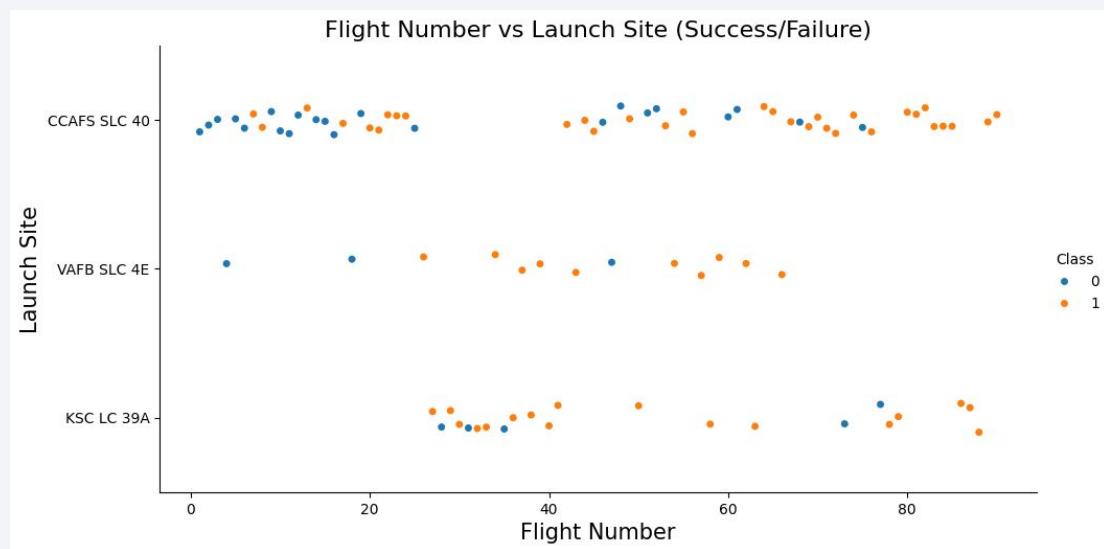
- Exploratory data analysis results
- Interactive analytics demo in screenshots
- Predictive analysis results

Section
2

Insights drawn from EDA

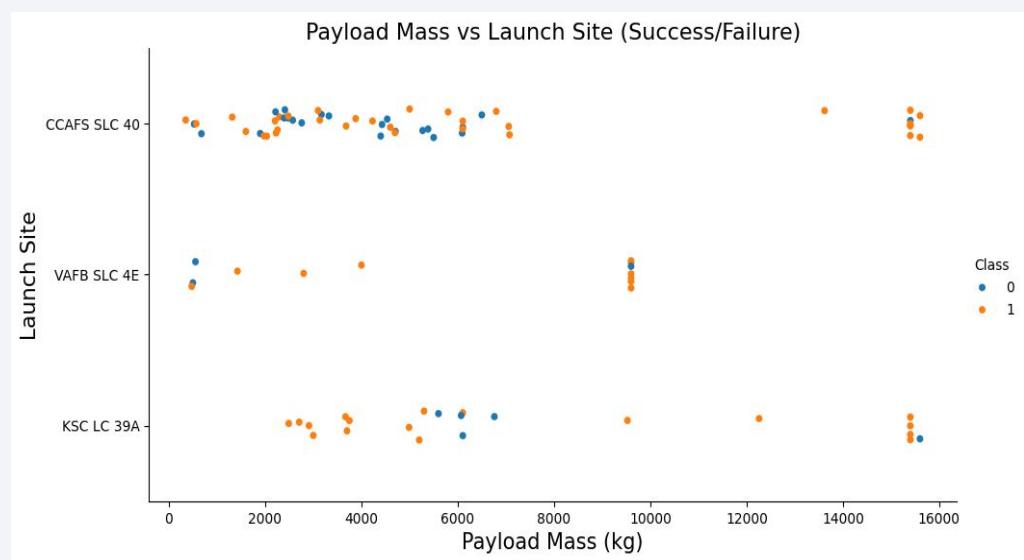
Flight Number vs. Launch Site

- The plot showed that launch sites with a higher number of flights tend to have higher success rates.



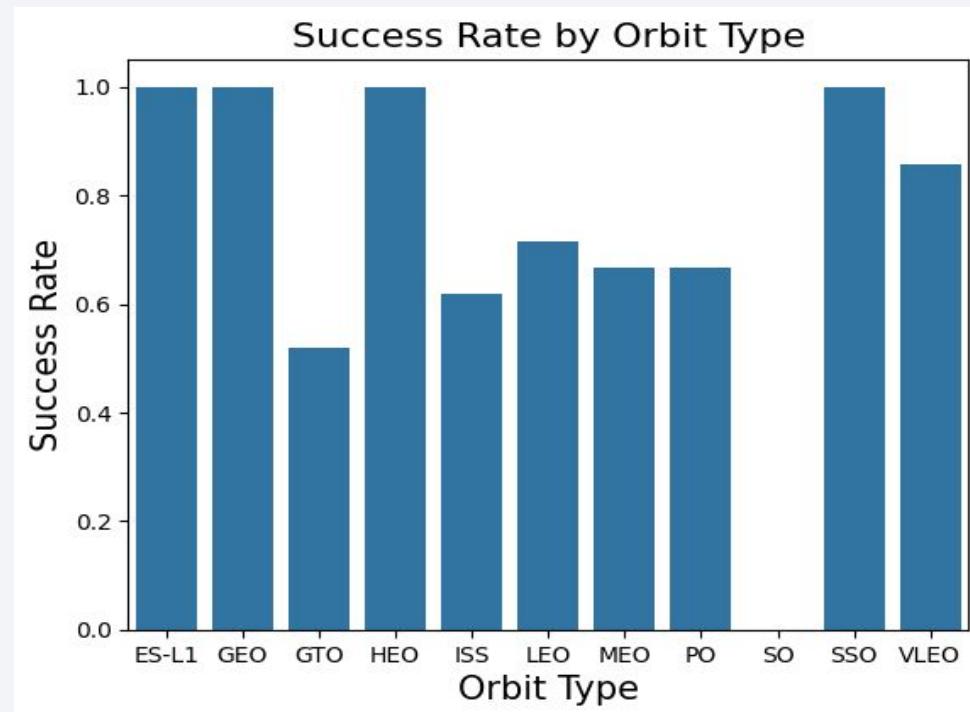
Payload vs. Launch Site

- CCAFS SLC 40 demonstrates the greatest payload versatility, successfully launching missions ranging from ~1,000 kg to over 15,000 kg, while VAFB SLC 4E and KSC LC 39A operate within more constrained payload mass ranges.



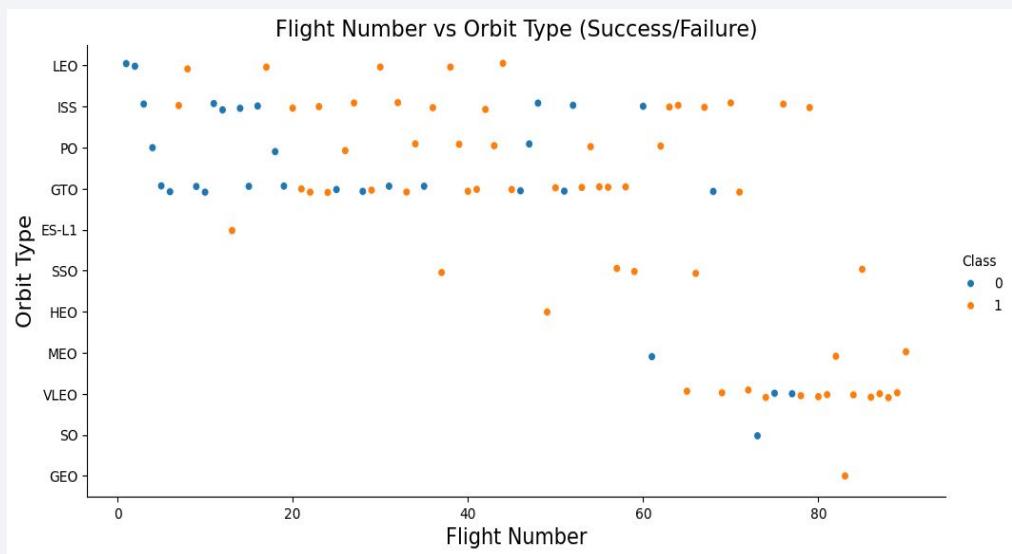
Success Rate vs. Orbit Type

- Mission success rates vary significantly by orbit type, with ES-L1, GEO, HEO, and SSO achieving near-perfect reliability (~100%) while GTO and ISS missions show notably lower success rates (~52-62%).



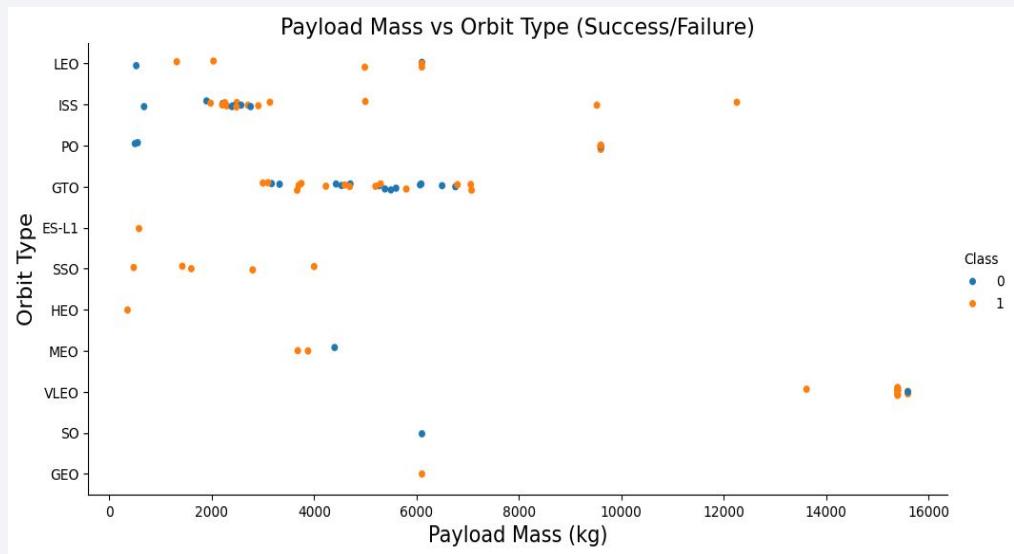
Flight Number vs. Orbit Type

- Mission complexity and frequency patterns emerge over time, with GTO and ISS missions dominating early flights, while specialized orbits like VLEO and SSO become more prominent in later flight numbers, suggesting operational maturity and expanded mission capabilities.



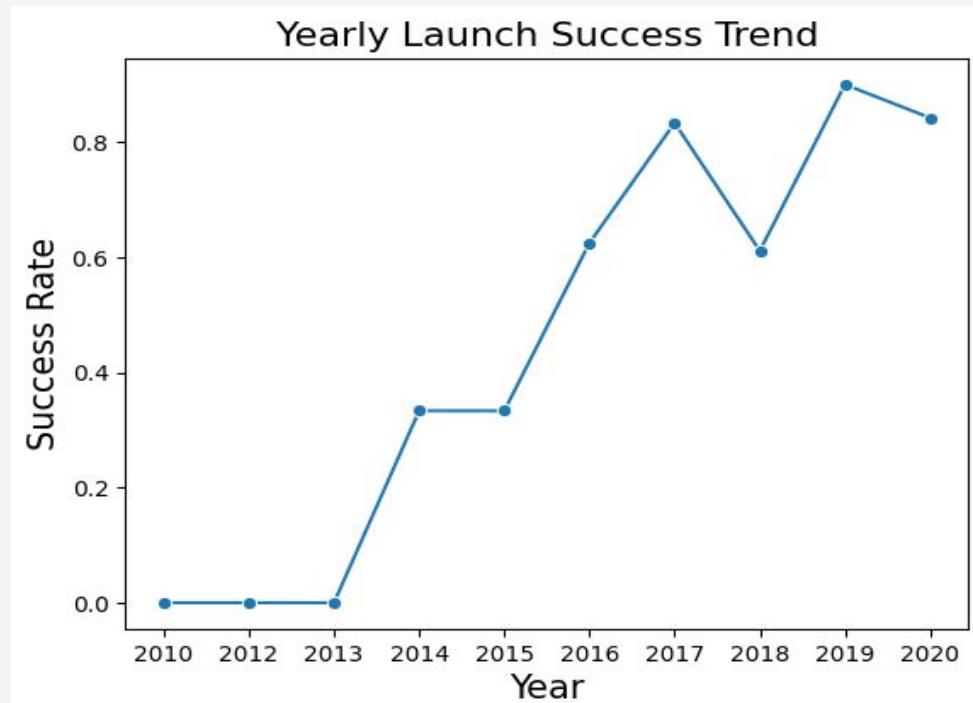
Payload vs. Orbit Type

- Different orbit types accommodate distinct payload mass ranges, with ISS missions clustering around 2,000-3,000 kg, GTO missions spanning 3,000-7,000 kg, and VLEO missions handling the heaviest payloads at 14,000-15,000 kg, reflecting specific mission requirements and orbital mechanics constraints.



Launch Success Yearly Trend

- Launch success rates demonstrate steady improvement over time, evolving from 0% in early years (2010-2013) to consistently high performance above 80% by 2019-2020, indicating significant operational learning and system maturation.



All Launch Site Names

- We used the keyword DISTINCT to filter out the unique launch site names.

Task 1

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

```
▷ %sql SELECT DISTINCT "Launch_Site" FROM SPACEXTABLE;
[10] ✓ 0.0s
... * sqlite:///my_data1.db
Done.

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

Python

Launch Site Names Begin with 'CCA'

- We used the LIMIT keyword set to 5 to retrieve the first five records of launch sites whose names start with 'CCA'.

Task 2

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

```
[11] %sql SELECT * FROM SPACEXTABLE WHERE "Launch_Site" LIKE 'CCA%' LIMIT 5;
[11] ✓ 0.0s
...
* 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

- Using the query below, we determined that the total payload carried by NASA boosters is 45,596.

Task 3

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

+ Code + Markdown

```
[12] %sql SELECT SUM("PAYLOAD_MASS_KG_") AS Total_Payload_Mass FROM SPACEXTABLE WHERE "Customer" = 'NASA (CRS)';
```

Python

✓ 0.0s

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

```
... Total_Payload_Mass
```

45596

Average Payload Mass by F9 v1.1

- We found that the average payload mass for the F9 v1.1 booster version is 2,534.67 kg.

Task 4

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

```
▷ %sql SELECT AVG("PAYLOAD_MASS__KG_") AS Average_Payload_Mass FROM SPACEXTABLE WHERE "Booster_Version" LIKE 'F9 v1.1%';
[13] ✓ 0.0s
... * sqlite:///my_data1.db
Done.

... Average_Payload_Mass
2534.6666666666665
```

Python

[+ Code](#) [+ Markdown](#)

First Successful Ground Landing Date

- The first successful ground landing date was: 22 December 2015.

Task 5

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

Hint:Use min function

```
%sql SELECT MIN(Date) AS First_Success_Ground_Pad FROM SPACEXTABLE WHERE "Landing_Outcome" = 'Success (ground pad)';
```

[14] ✓ 0.0s

Python

```
... * sqlite:///my\_data1.db
```

Done.

```
... First_Success_Ground_Pad
```

2015-12-22

Successful Drone Ship Landing with Payload between 4000 and 6000

- We applied the WHERE clause to filter boosters that successfully landed on a drone ship, and used the AND condition to select those with a payload mass between 4,000 and 6,000 kg.

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 SPACEXTABLE WHERE "Landing_Outcome" = 'Success (drone ship)' AND "PAYLOAD_MASS_KG_" > 4000 AND "PAYLOAD_MASS_KG_" < 6000;
```

[15] ✓ 0.0s

* [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

- This query counts the number of occurrences for each unique mission outcome, providing a summary of how many launches resulted in success, failure, or other outcomes.

Task 7

List the total number of successful and failure mission outcomes

```
%sql SELECT "Mission_Outcome", COUNT(*) AS Outcome_Count FROM SPACETABLE GROUP BY "Mission_Outcome";  
[16] ✓ 0.0s  
... * sqlite:///my\_data1.db  
Done.  
...  


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


```

Boosters Carried Maximum Payload

- We identified the booster that carried the maximum payload by using the MAX() function along with a subquery in the WHERE clause.

Task 8

List all the booster_versions that have carried the maximum payload mass. Use a subquery.

```
▶ %sql SELECT "Booster_Version" FROM SPACEXTABLE WHERE "PAYLOAD_MASS_KG_" = (SELECT MAX("PAYLOAD_MASS_KG_") FROM SPACEXTABLE);
[17] ✓ 0.0s
...
* 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
```

Python

2015 Launch Records

- We combined the WHERE clause with LIKE, AND, and BETWEEN conditions to filter records of failed drone ship landings in 2015, retrieving the corresponding booster versions and launch site names.

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.

```
[18]: %sql SELECT substr(Date, 6, 2) AS Month, "Landing_Outcome", "Booster_Version", "Launch_Site" FROM SPACEXTABLE WHERE "Landing_Outcome" = 'Failure (drone ship)' AND substr(Date, 0, 5) = '2015';
      ✓ 0.0s
...
* sqlite:///my_data1.db
Done.

...
   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
```

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

- We queried the landing outcomes and their counts, filtering records between 2010-03-20 and 2010-06-04 using the WHERE clause. The results were grouped by landing outcome using GROUP BY and sorted in descending order with the ORDER BY clause.

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 Outcome_Count FROM SPACETABLE WHERE Date BETWEEN '2010-06-04' AND '2017-03-20' GROUP BY "Landing_Outcome" ORDER BY Outcome_Count DESC;
[19] ✓ 0.0s
... * sqlite:///my\_data1.db
Done.

...

| 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             |
| Precidited (drone ship) | 1             |


```

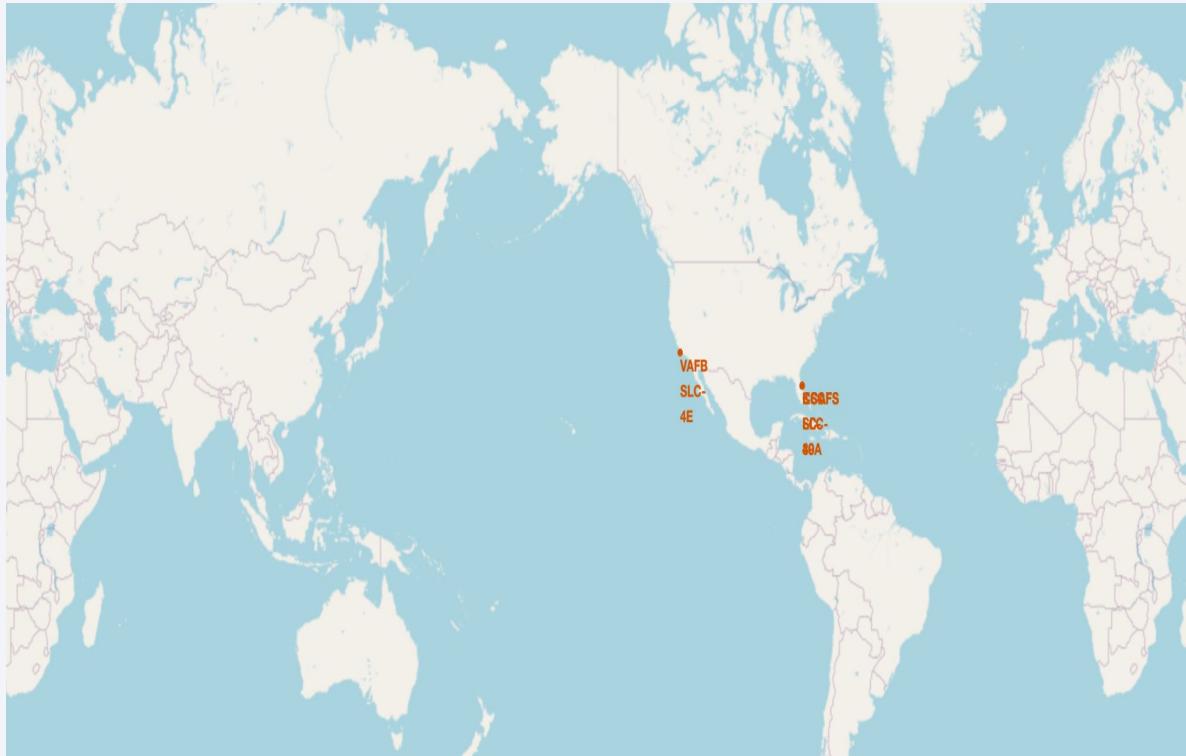
Python

The background of the slide is a photograph taken from space at night. It shows the curvature of the Earth's horizon against a dark blue sky. City lights are visible as numerous small white and yellow dots, primarily concentrated in the lower right quadrant where the United States appears. In the upper right, there is a bright green and yellow glow, likely representing the Aurora Borealis or a similar atmospheric phenomenon.

Section
3

Launch Sites Proximities Analysis

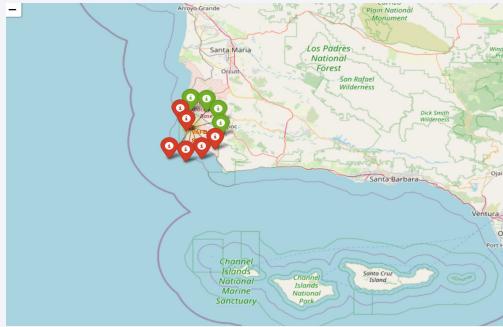
All Launch Sites global map markers



- The launch sites for spacex are along the borders.

Map showing launch sites with successes and failures

- The green markers show success and the red ones show failure.



Launch Site distance to Landmarks

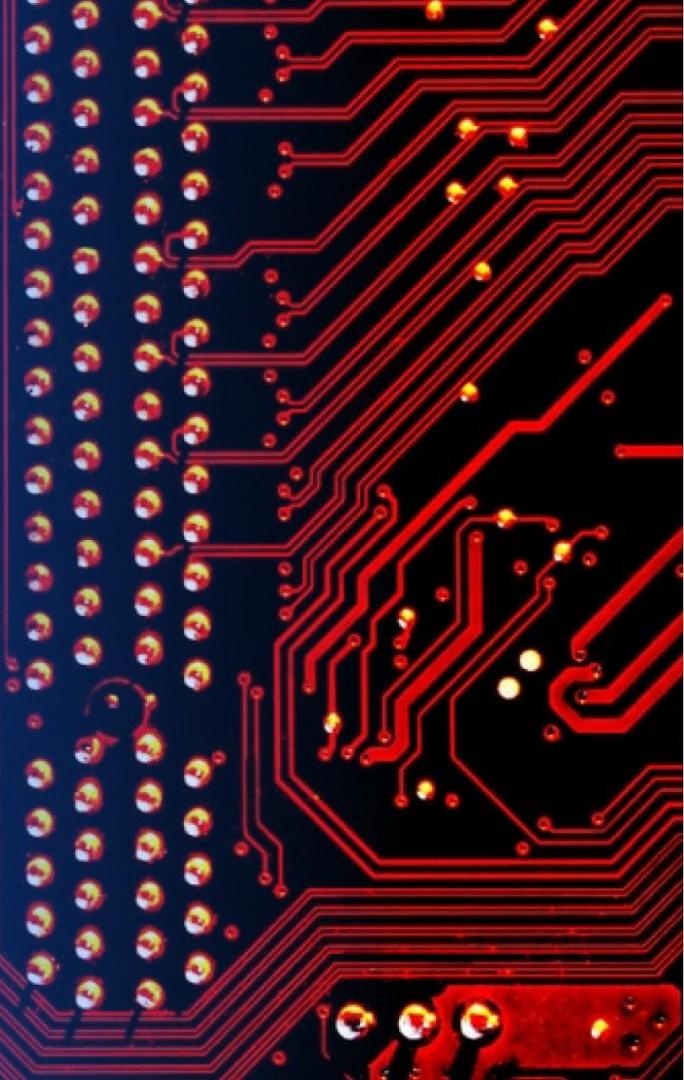


The launch sites are near the coast and are away from major cities. They are also not in close proximities to the highway or railways.

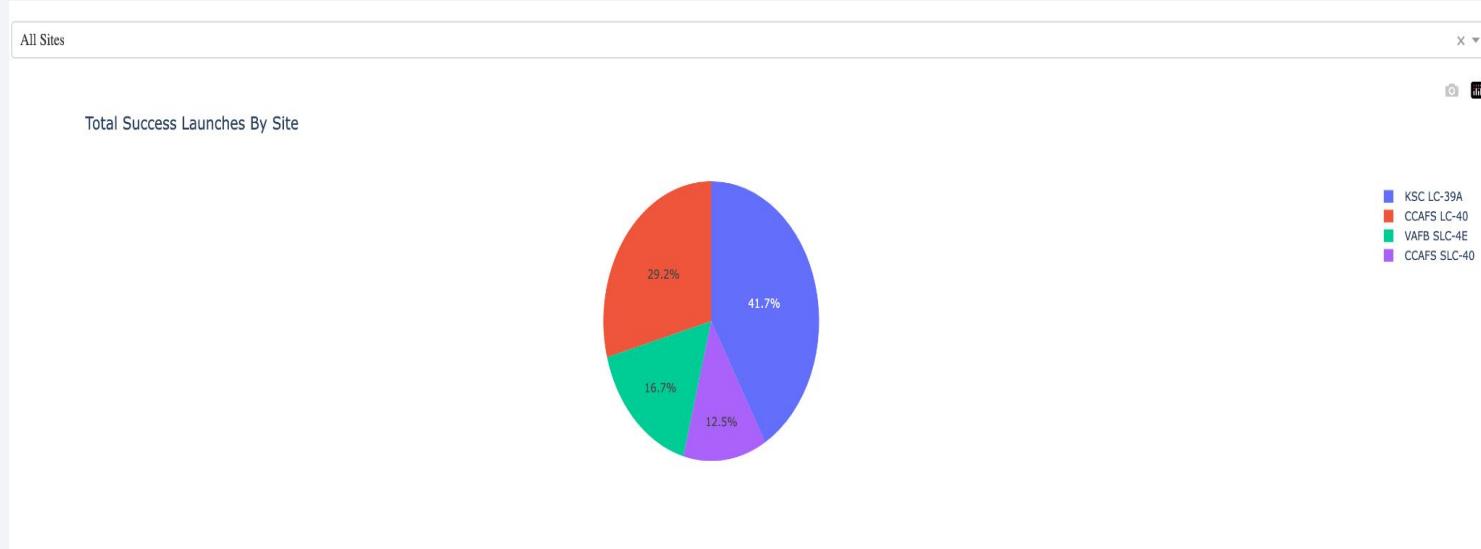
Coast: 0.9km; City: 78.45km; highway: 29.21km

Section
4

Build a Dashboard with Plotly Dash

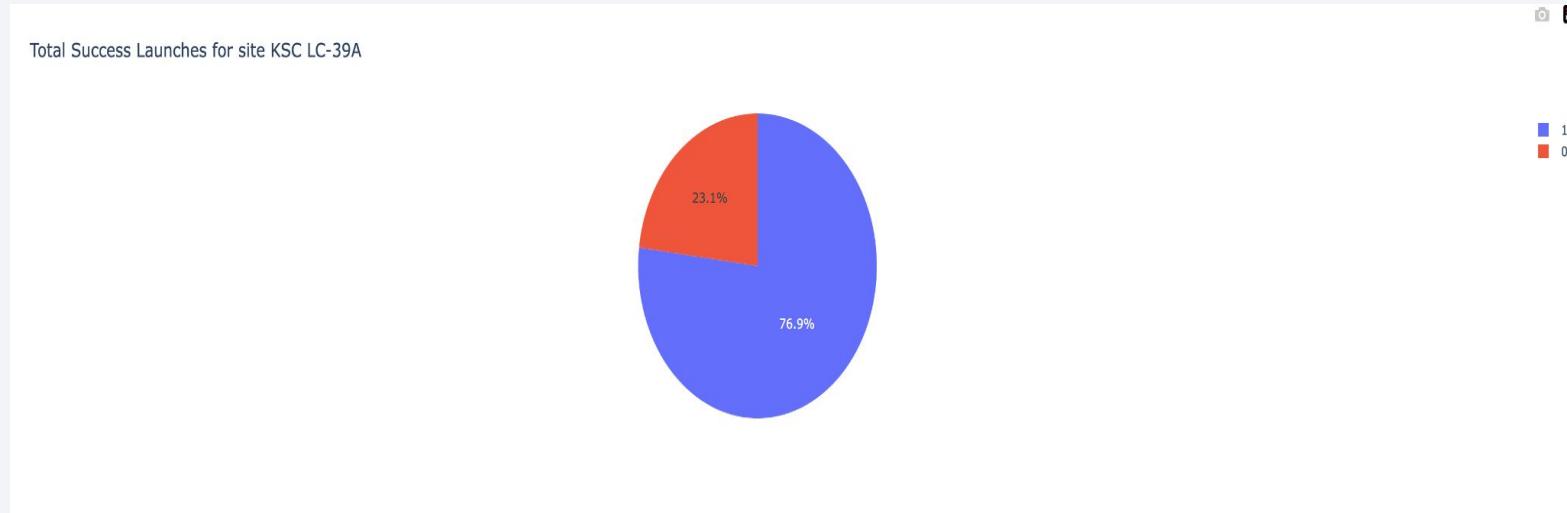


Pie chart showing the success percentage achieved by each launch site



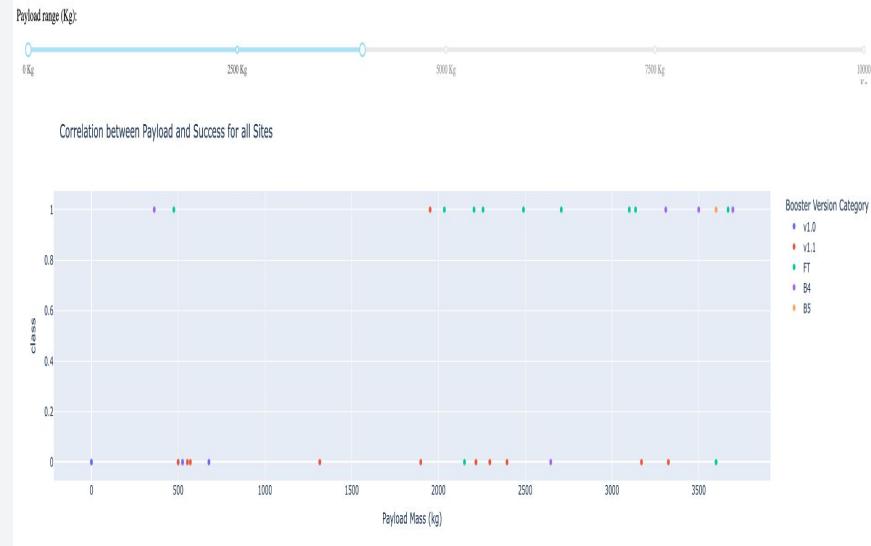
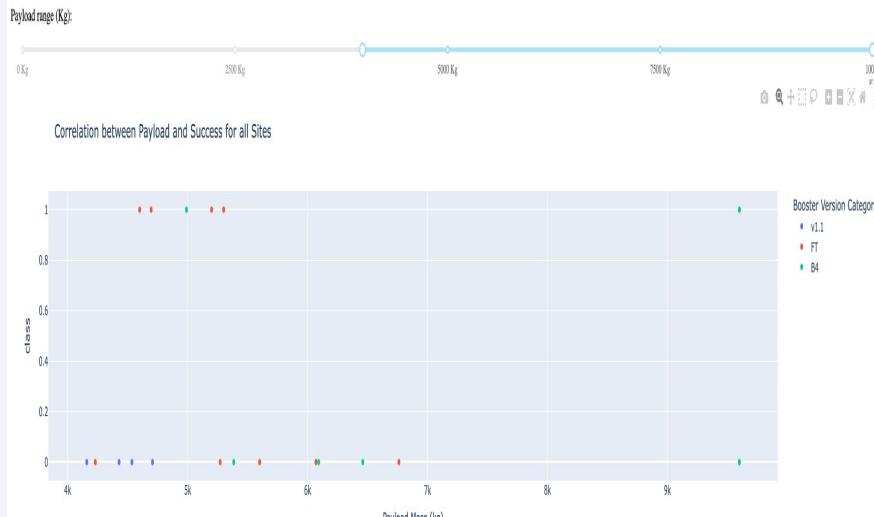
It can be noted that launch site kSC LC-39A had the highest success percentage.

Pie chart showing the Launch site with the highest launch success ratio



The KSC LC-39A achieved a success percentage is 76.9% while the failure is 23.1%.

Scatter plot of Payload vs Launch Outcome for all sites



Success for low weighted payloads is greater than that of heavier payloads

Section
5

Predictive Analysis (Classification)

Classification Accuracy

- For the test set, we saw that Logistic Regression, SVM, and KNN all tie for the best accuracy at 83.33% while the decision tree had an accuracy of 77.78%.
- This maybe due to the small test set size. The scores of the entire data set shows that the decision tree is the best model. This model not only has high scores but high accuracy as well.

TASK 12

Find the method performs best:

```
# Print the test accuracies for all models
print("Logistic Regression Test Accuracy:", logreg_cv.score(X_test, Y_test))
print("SVM Test Accuracy:", svm_cv.score(X_test, Y_test))
print("Decision Tree Test Accuracy:", tree_cv.score(X_test, Y_test))
print("KNN Test Accuracy:", knn_cv.score(X_test, Y_test))

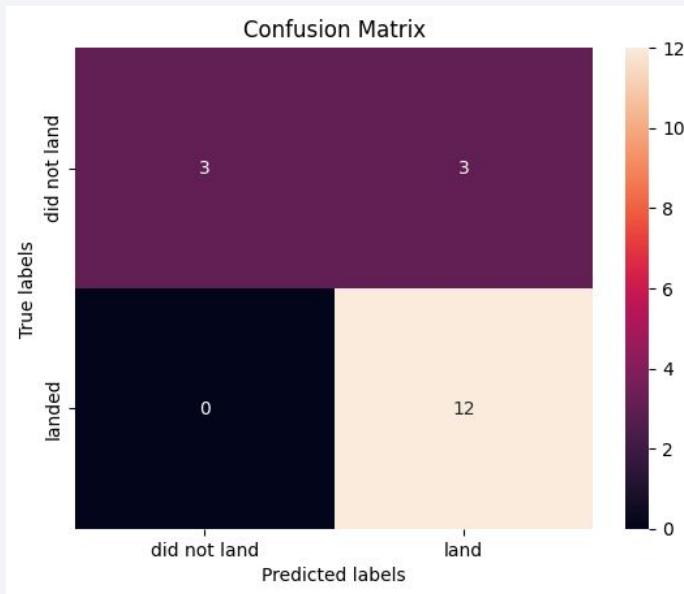
[31] ✓ 0.0s
...
Logistic Regression Test Accuracy: 0.8333333333333334
SVM Test Accuracy: 0.833333333333334
Decision Tree Test Accuracy: 0.7777777777777778
KNN Test Accuracy: 0.833333333333334
```

```
models = {'KNeighbours': knn_cv.best_score_, 'Decision Tree': tree_cv.best_score_, 'Logistic Regression': logreg_cv.best_score_, 'Support Vector Machine': svm_cv.best_score_}
best_model = max(models, key=models.get)
print('Best model is', best_model, 'with a score of', models[best_model])
if best_model == 'KNeighbours':
    print('KNeighbours is the best model with params', knn_cv.best_params_)
elif best_model == 'Decision Tree':
    print('Decision Tree is the best model with params', tree_cv.best_params_)
elif best_model == 'Logistic Regression':
    print('Logistic Regression is the best model with params', logreg_cv.best_params_)
elif best_model == 'Support Vector Machine':
    print('Support Vector Machine is the best model with params', svm_cv.best_params_)

[33] ✓ 0.0s
...
Best model is Decision Tree with a score of 0.875
Decision Tree is the best model with params {'criterion': 'gini', 'max_depth': 6, 'max_features': 'sqrt', 'min_samples_leaf': 1, 'min_samples_split': 10, 'splitter': 'random'}
```

Confusion Matrix

- The identical performance and matching confusion matrices of the three models may be due to the limited size of the training dataset, which contained only 18 samples, making it difficult to distinguish their effectiveness.



Conclusions

We can conclude the following:

- Launch sites with a higher number of flights tend to have higher success rates.
- Launch success rates began improving steadily from 2013 through 2020.
- Orbits such as ES-L1, GEO, HEO, SSO, and VLEO showed the highest success rates.
- KSC LC-39A recorded the most successful launches among all sites.
- Decision Tree Model is the best algorithm for this dataset.

Thank you!

