

Winning Space Race with Data Science

G. Marchal
27/05/2023



Outline

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

Executive Summary

We here studied SpaceX Falcon 9 launch data in order to predict whether the rocket's first stage will land successfully.

Several Machine Learning classification algorithms were used, following the following steps:

- Data Collection, Wrangling, and Formatting
- Exploratory Data Analysis
- Interactive Data Visualisation
- Predictive Machine Learning

Main findings:

- Many features of the Falcon 9 Launch influence the outcome of the landing

Introduction

- Project background and context
 - Space industry is becoming more mainstream, thanks to private companies entering the market. Analysis of launch costs and parameters influencing mission success are vital to keep this market viable
 - SpaceX developed a rocket with an unique re-usable first stage (RFS), cutting the costs of launches by more than 100 million USD when compared to competitors
 - Optimising the success rate of RFS preservation after launches is essential to maintain this financial edge
- Problems you want to find answers
 - Determine successful landing of the RFS
 - Assess the impact of various factors on landing outcome
 - Analyse correlations between launch sites and landing outcome

Section 1

Methodology

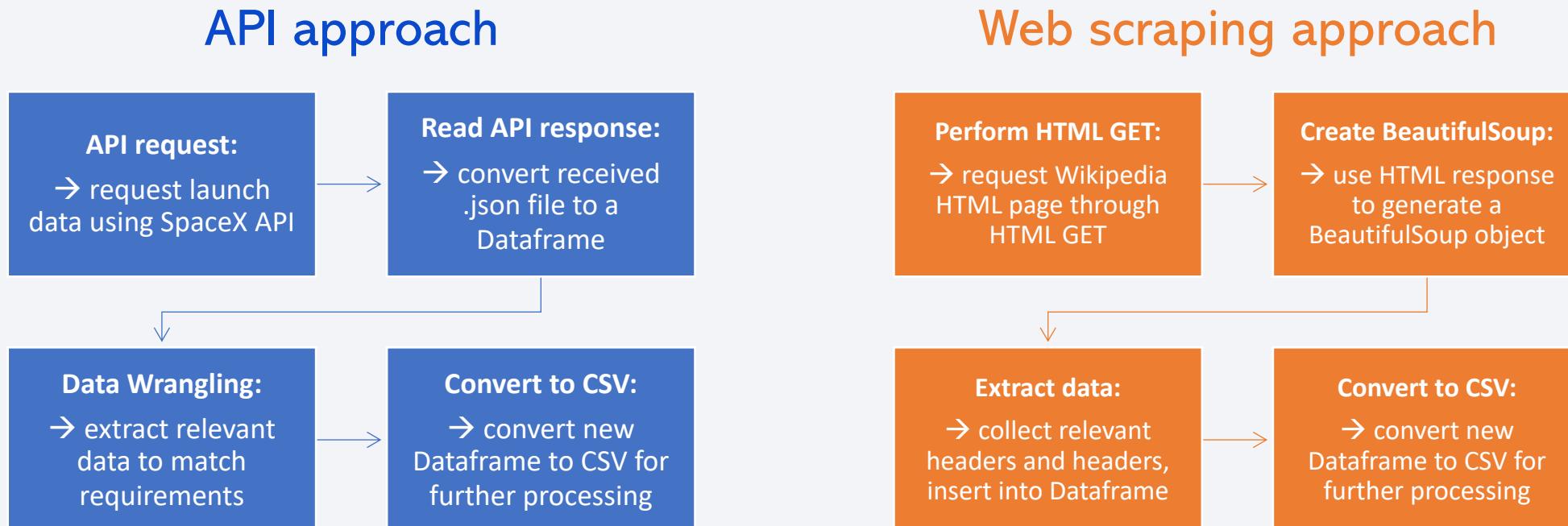
Methodology

Executive Summary

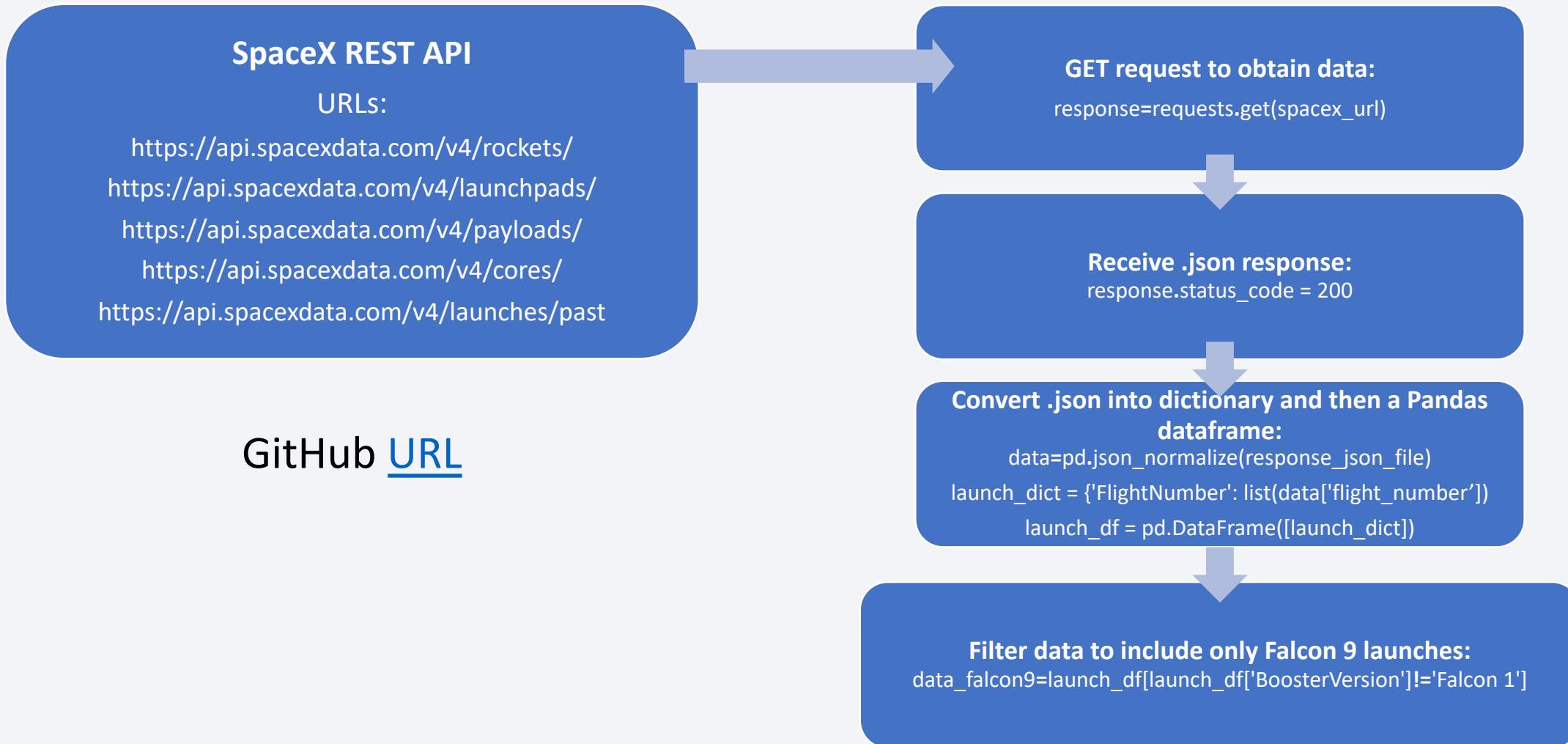
- Data collection methodology:
 - SpaceX API
 - Web scraping of launch records from relevant Wikipedia pages
- Perform data wrangling
 - Supervised models were trained after converting mission outcomes (0: unsuccessful, 1: successful)
- Perform exploratory data analysis (EDA) using visualisation and SQL
- Perform interactive visual analytics using Folium and Plotly Dash
- Perform predictive analysis using classification models
 - Created a column for 'class'
 - Standardised and transformed data
 - Split dataset in train/test data
 - Determination of best classification algorithm using test data

Data Collection

- Data was collected using the SpaceX API and by web scraping the Wikipedia page listing all Falcon 9 (heavy) launches



Data Collection – SpaceX API



Data Collection - Scraping

Wikipedia page

"List of Falcon 9 and Falcon Heavy launches"

```
static_url =
```

```
https://en.wikipedia.org/w/index.php?title=List_of_Falcon_9_and_Falcon_Heavy_launches&oldid=1027686922
```

```
element = soup.find_all('th')
for row in range(len(element)):
    try:
        name = extract_column_from_header(element[row])
        if (name is not None and len(name) > 0):
            column_names.append(name)
    except:
        pass
```

GET request to obtain data:

```
html_data = requests.get(static_url)
```

Create object from response with BeautifulSoup:

```
soup = BeautifulSoup(html_data.text, 'html.parser')
```

Find HTML table with Falcon 9 data:

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

Iterate through <th> elements to extract column names

Create and insert data in Pandas dataframe

```
launch_dict= dict.fromkeys(column_names)

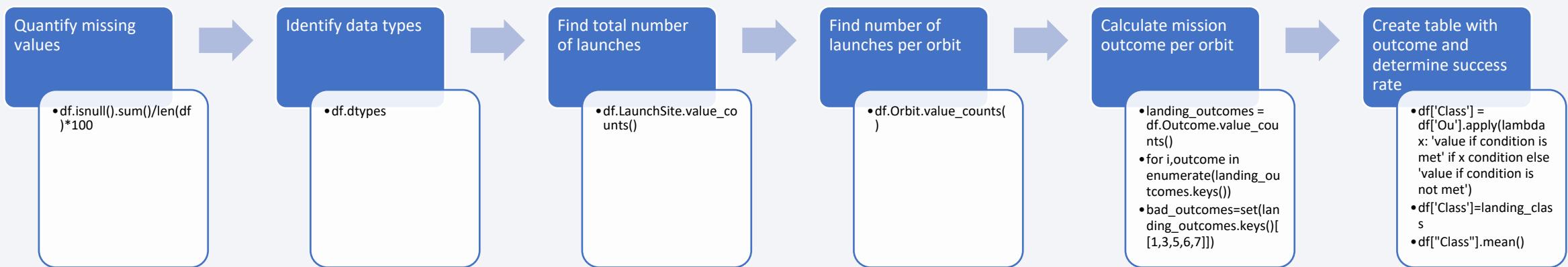
# Remove an irrelevant column
del launch_dict['Date and time ( )']

# Let's initial the launch_dict with each
launch_dict['Flight No.']= []
launch_dict['Launch site']= []
launch_dict['Payload']= []
launch_dict['Payload mass']= []
launch_dict['Orbit']= []
launch_dict['Customer']= []
launch_dict['Launch outcome']= []
# Added some new columns
launch_dict['Version Booster']= []
launch_dict['Booster landing']= []
launch_dict['Date']= []
launch_dict['Time']= []
```

```
df= pd.DataFrame({ key:pd.Series(value) for key, value in launch_dict.items() })
```

Data Wrangling

- Data was first assessed by identifying missing points, data type, and count
- Then, data was split per orbit type and mission outcome per orbit was calculated
- Data was subsequently put into a table, and success rate per orbit type was calculates



EDA with Data Visualization

Plot types used:

- Scatter plots:
 - To identify relationships between two variables
 - Flight Number vs. Payload
 - Flight Number vs. Launch Site
 - Payload vs. Launch Site
 - Class vs. Orbit
 - Flight Number vs. Orbit
 - Payload vs. Orbit
- Bar charts:
 - To compare values between two groups, often used to compare a variable at a given point in time
 - Success Rate per Orbit
- Line charts:
 - To track changes over time
 - Success Rate over Time

EDA with SQL

SQL queries performed to explore the data:

- Identify unique Launch Sites:
 - %sql `SELECT DISTINCT "Launch_Site" FROM SPACEXTABLE;`
- Retrieve 5 records with Launch Site name beginning with CCA:
 - %sql `SELECT "Launch_Site" FROM SPACEXTABLE WHERE "Launch_Site" LIKE "CCA%" LIMIT 5;`
- Find total Payload Mass carried by NASA-launched boosters (CRS):
 - %sql `SELECT customer, SUM(PAYLOAD_MASS_KG_) AS total_payload FROM SPACEXTABLE WHERE customer="NASA (CRS)" GROUP BY customer;`
- Retrieve average Payload Mass carried by F9 v1.1 boosters
 - %sql `SELECT "Booster_Version", AVG(PAYLOAD_MASS_KG_) AS avg_payload FROM SPACEXTABLE WHERE "Booster_Version"="F9 v1.1" GROUP BY "Booster_Version";`
- List date with first successful Landing Outcome:
 - %sql `SELECT "Landing_Outcome", MIN("Date") AS min_date FROM SPACEXTABLE WHERE "Landing_Outcome"="Success (ground pad)" GROUP BY "Landing_Outcome";`
- List Boosters with success in Drone Ship and a Payload Mass >4000 and <6000:
 - %sql `SELECT DISTINCT "Booster_Version" FROM SPACEXTABLE WHERE "Landing_Outcome"="Success (drone ship)" AND (PAYLOAD_MASS_KG_>4000 AND PAYLOAD_MASS_KG_<6000);`
- List total number of successful and failure Mission Outcomes:
 - %sql `SELECT "Mission_Outcome", COUNT(*) AS count_mission_outcome FROM SPACEXTABLE GROUP BY TRIM("Mission_Outcome");`
- Find names of Booster Versions which carried the max. Payload Mass (using a subquery):
 - %sql `SELECT "Booster_Version" FROM SPACEXTABLE WHERE PAYLOAD_MASS_KG_ = (SELECT MAX(PAYLOAD_MASS_KG_) AS max_payload FROM SPACEXTABLE);`
- Retrieve records, displaying Month Name, Landing Outcome, and Launch Site for missions in 2015:
 - %sql `SELECT substr("Date",1,4) AS year, 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';`
- Rank the count of Landing Outcomes between 2010-06-04 and 2017-03-20:
 - %sql `SELECT Count("Landing_Outcome") AS count_landing_outcomes FROM SPACEXTABLE WHERE ("Landing_Outcome" = "Failure (drone ship)" OR "Landing_Outcome"="Success (ground pad)") AND ("Date" BETWEEN "2010-06-04" AND "2017-03-20")`

Build an Interactive Map with Folium

Folium is a Python library allowing creation of interactive maps

- We created a Folium map containing:
 - Circles to highlight Launch Sites
 - Markers to indicate Mission Outcome
 - Indicator of Mouse Pointer position
- Additionally, we used lines and markers to calculate the distance of Launch Sites and:
 - Railways
 - Highways
 - Coastlines
 - Cities

Build a Dashboard with Plotly Dash

We built a Plotly Dash dashboard to visualise data in a real-time, interactive manner

Visualisations used:

- Pie chart to visualise contributions
 - Total Mission Success Counts per Launch Site
 - Identify Launch Site with highest Success Ratio
- Scatter plots to visualise relationships between variables:
 - Mission Outcome vs Payload Mass

Predictive Analysis (Classification)

Predictive analysis was used to predict the Mission Outcome

Dependent Variable:

- Class (Mission Outcome; 1 = success, 0 = fail)

Excluded from Independent Variables:

- Date, Outcome, Booster Version, Longitude, Latitude

One-hot encoding of categorical Independent Variables:

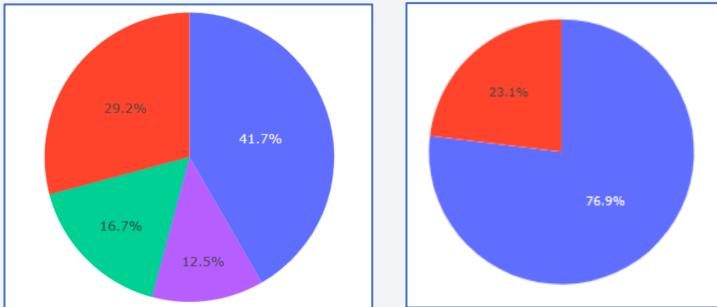
- Orbit, Launch Site, Landing Pad, Serial. Grid Fins, Reused, Legs

- Data was standardised, split into train and test data, and various classification methods were assessed.
- Optimal Hyperparameters per Model and the best performing Model were identified



Results

- Exploratory data analysis results
 - Relevant parameters influencing Mission Outcome were identified
- Interactive analytics demo screenshots



- Predictive analysis results

Model Name	Score
0 Logistic Regression	0.833333
1 SVM	0.833333
2 Decision Tree	0.833333
3 K Nearest Neighbors	0.833333

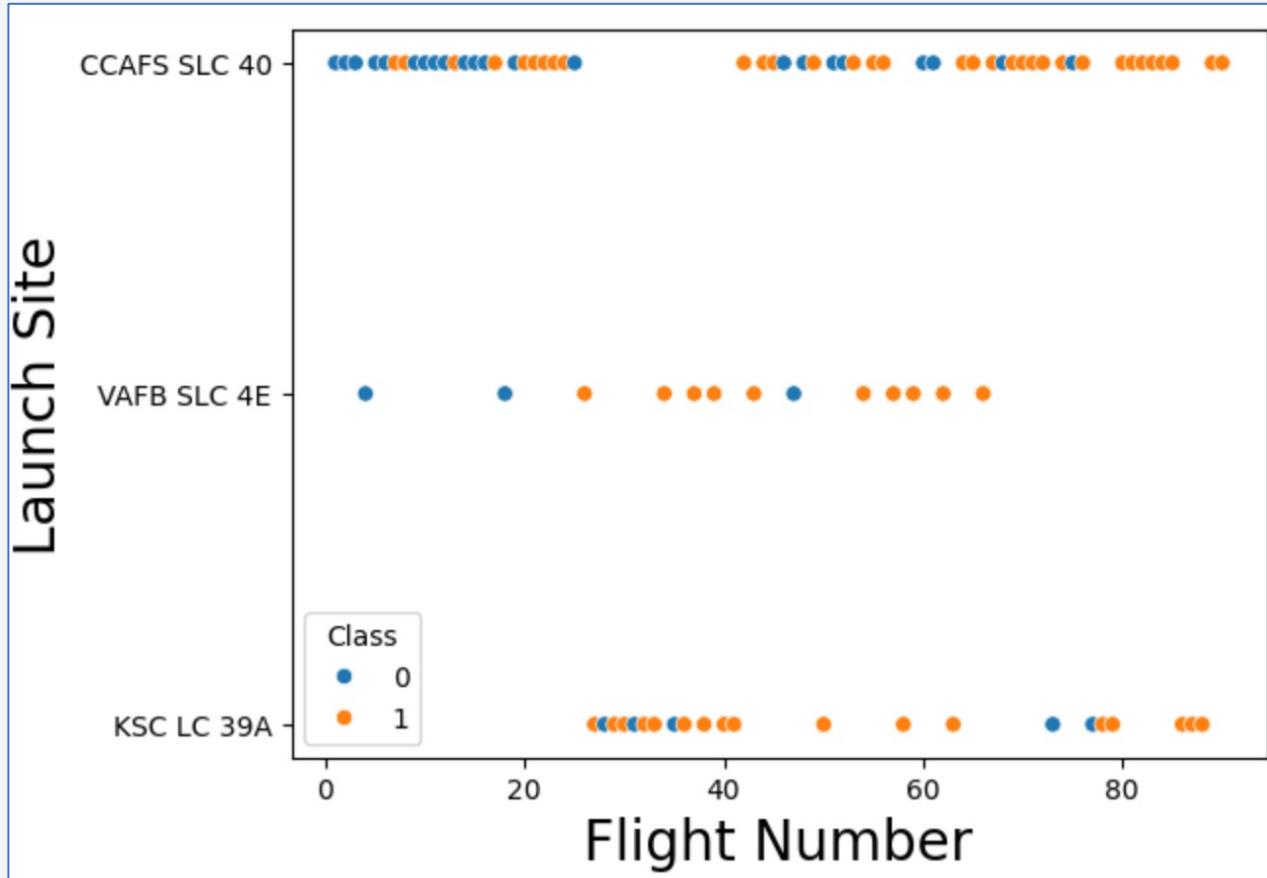
All models have the same score of 0.833. The confusion matrix results are also the same.

The background of the slide features a complex, abstract digital visualization. It consists of numerous thin, glowing lines that create a sense of depth and motion. The lines are primarily blue and red, with some green and purple highlights. They form a grid-like structure that curves and twists across the frame, resembling a three-dimensional space or a network of data points. The overall effect is futuristic and dynamic.

Section 2

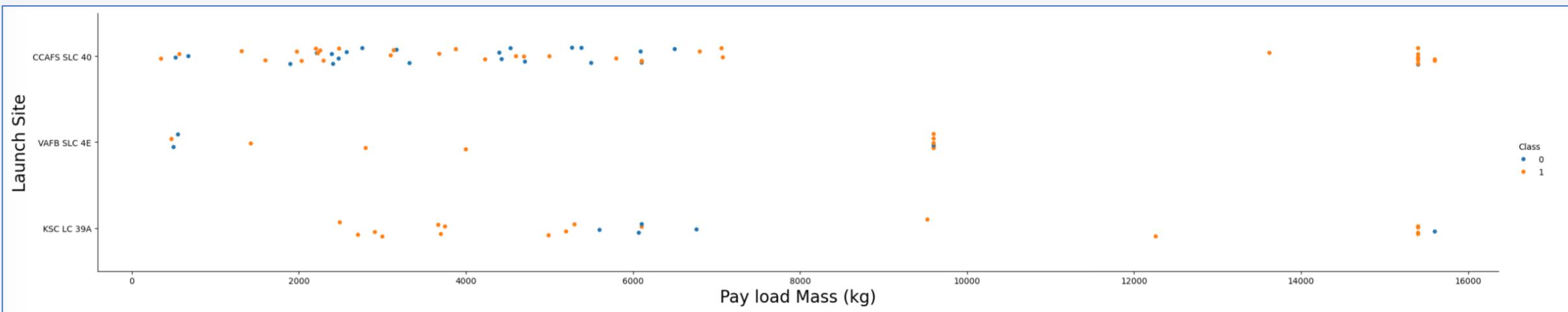
Insights drawn from EDA

Flight Number vs. Launch Site



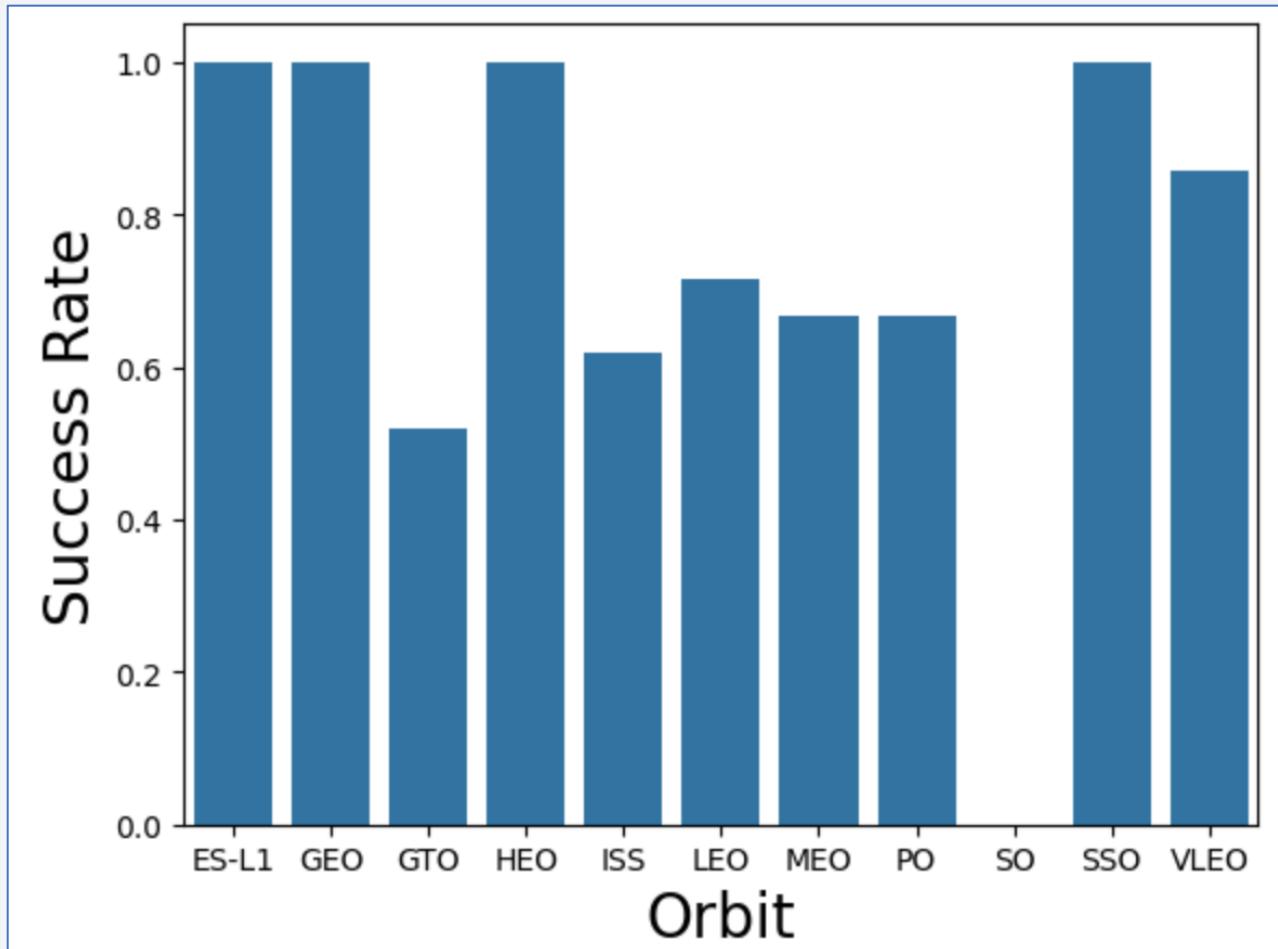
- Most initial missions were unsuccessful (coloured blue)
- As flight number increases, the mission is more likely to be successful (coloured orange)

Payload vs. Launch Site



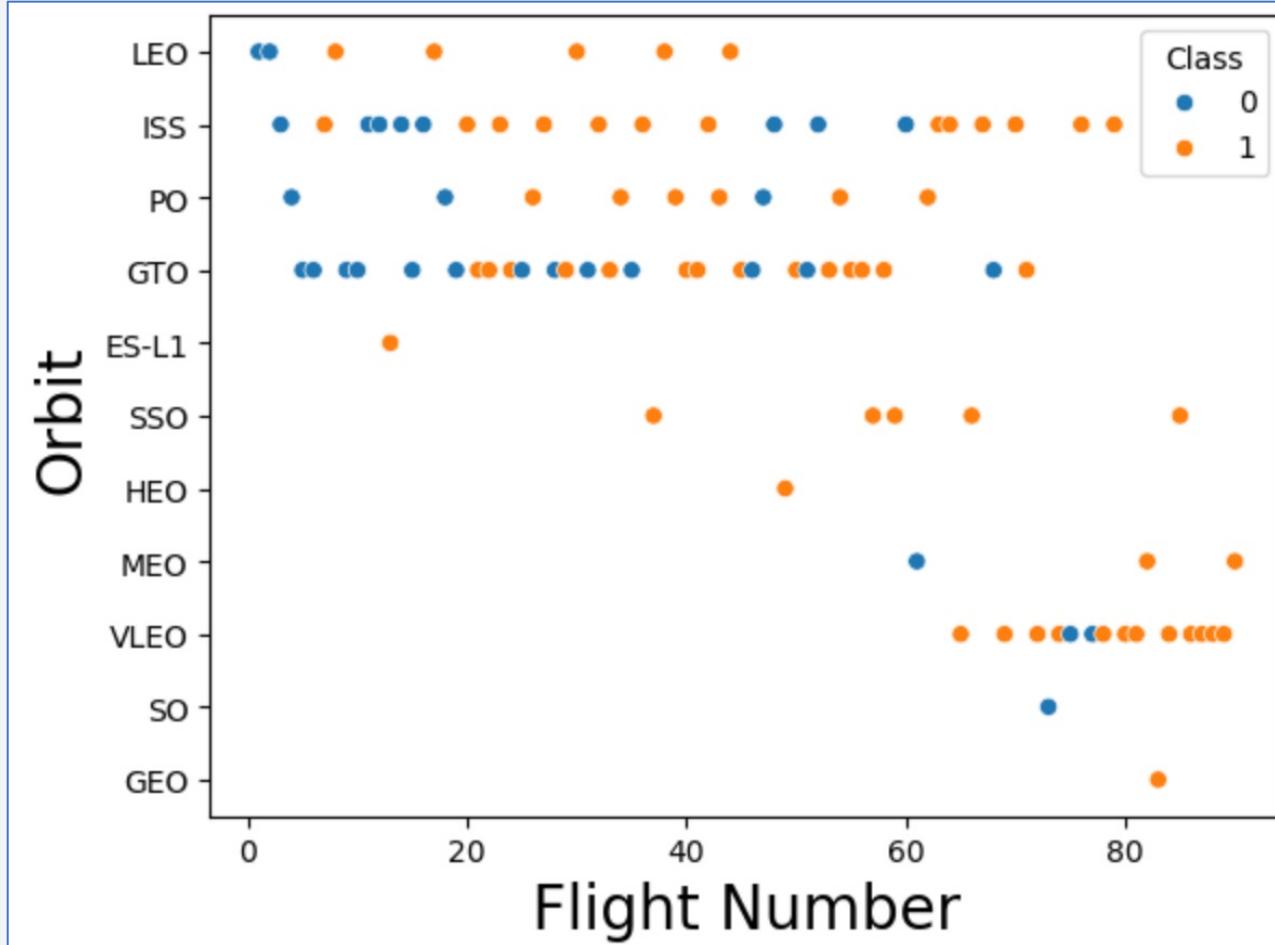
- At Launch Site 'CCAFS-SLC' (Caraveral Center), results for lighter Payloads (<8,000 kg) are mixed, but heavier Payloads are generally successful
- At Launch Site 'VAFB-SLC' (Vandenberg), the maximum Payload launched was 10,000 kg, and Missions were mostly successful
- At Launch Site 'KSC-LC' (Kennedy Space Center), most Missions are successful, apart from Missions with a Payload of around 6,000 kg

Success Rate vs. Orbit Type



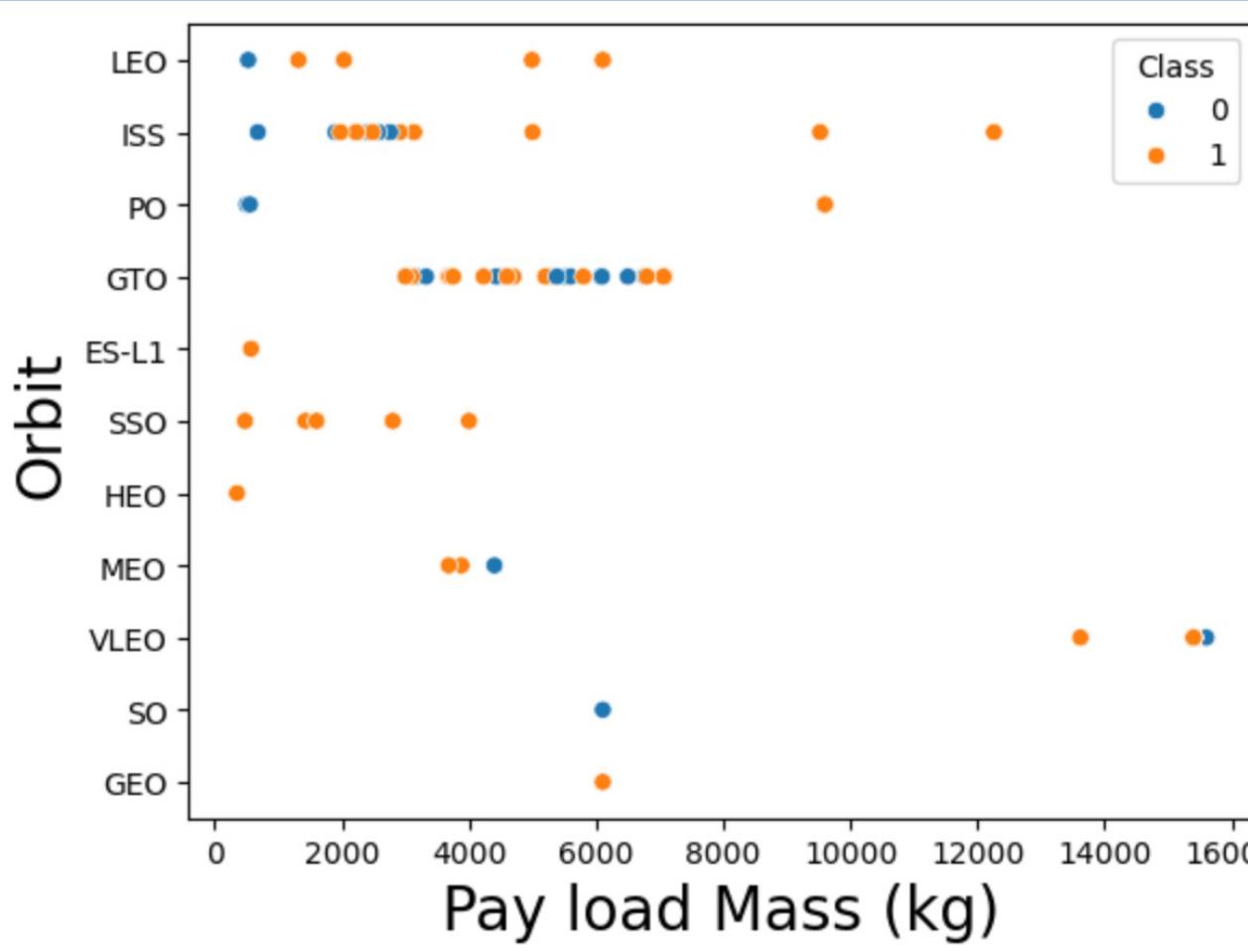
- Orbit types ES-L1, GEO, HEO, and SSO have a perfect success rate of 1
- VLEO also has a good success rate
- GTO, ISS, LEO, MEO, and PO have a mixed success rate
- All missions to orbit SO failed

Flight Number vs. Orbit Type



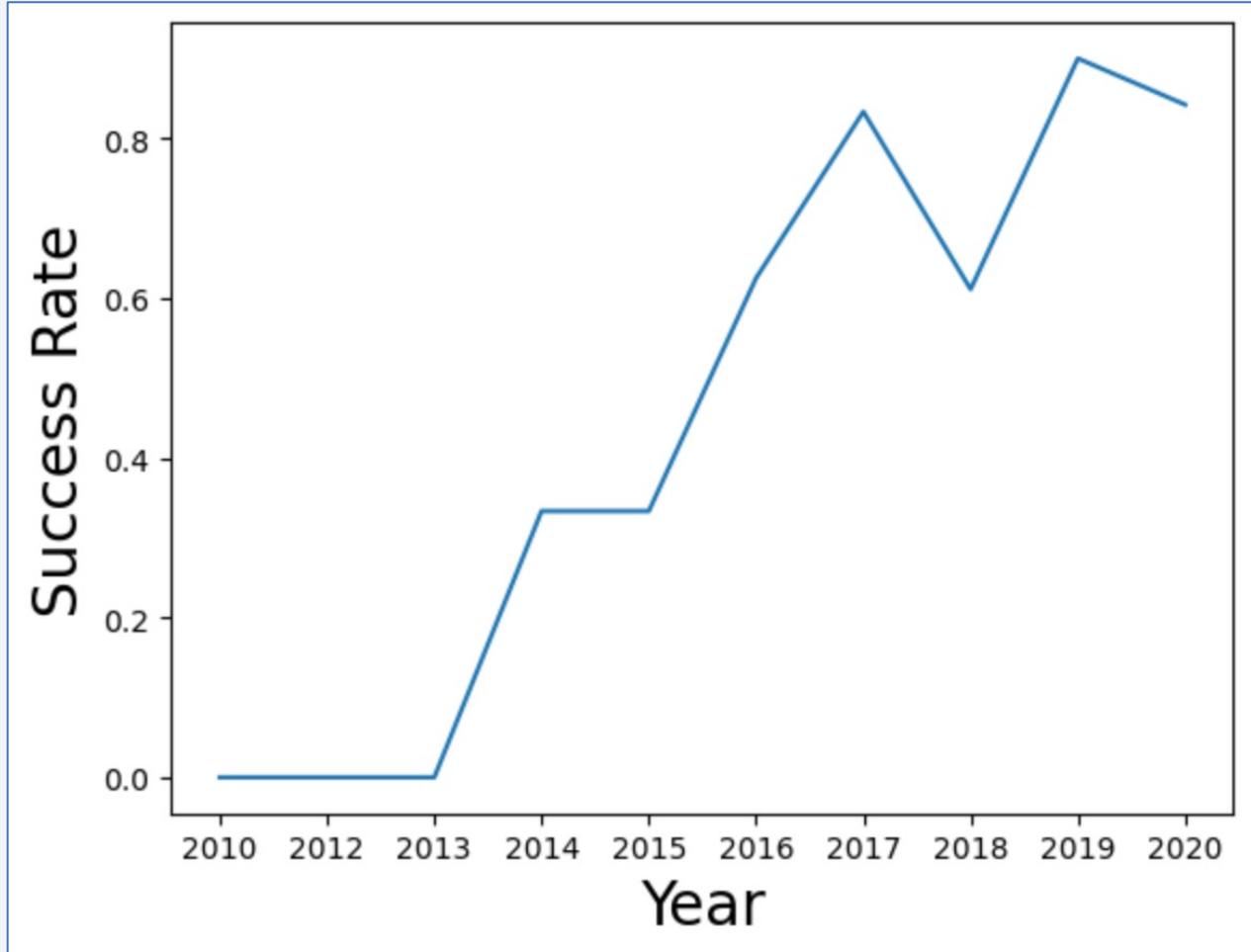
- Overall, flights with a low Flight Number are more often unsuccessful (coloured blue)
- After two failed Missions, all following Missions to LEO were successful (coloured orange)
- ISS displays a cluster of failed Missions around Flight Number 50
- For GTO, there is no clear relationship between Flight Number and Mission Outcome
- The 0% Success Rate for SO is due to only one Mission being attempted

Payload vs. Orbit Type



- Heavy Payloads ($>8,000$ kg) were only deployed to orbits ISS, PO, and VLEO.
- Missions with a heavy Payload were generally successful (coloured orange)
- Missions with a very light Payload (<1000 kg) are likely to be unsuccessful (coloured blue)
 - This is especially true for LEO, ISS, and PO
 - ES-L1, SSO, and HEO appear more suitable for very light Payloads
- For orbit GTO, there is no clear relationship between Payload and Mission Outcome

Launch Success Yearly Trend



- All Missions until 2013 were unsuccessful
- Mission success rate was the same in 2013 and 2015
- After 2015, Mission success rate started increasing, with the exemption of 2018

All Launch Site Names

```
%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 |


```

“DISTINCT” only returns unique values

Four Unique Launch Sites:

- CCAFS LC: Cape Canaveral Launch Complex
- VAFB SLC: Vandenberg Space Force Base
- KSC LC: Kennedy Space Center
- CCAFS SLC: Cape Canaveral Space Launch Complex

Launch Site Names Begin with 'CCA'

```
%sql SELECT "Launch_Site" FROM SPACEXTABLE WHERE "Launch_Site" LIKE "CCA%" LIMIT 5;  
* sqlite:///my_data1.db  
Done.  
Launch_Site  
_____  
CCAFS LC-40  
CCAFS LC-40  
CCAFS LC-40  
CCAFS LC-40  
CCAFS LC-40
```

Explanation:

*WHERE “Launch_Site” LIKE “CCA%” – to return Launch Sites starting with ‘CCA’
LIMIT 5 – to return only five records*

All returned values are CCAFS LC: Cape Canaveral Launch Complex

Total Payload Mass

```
%sql SELECT customer, SUM(PAYLOAD_MASS__KG_) AS total_payload FROM SPACEXTABLE WHERE customer="NASA (CRS)" GROUP BY customer;
```

* sqlite:///my_data1.db
Done.

Customer	total_payload
NASA (CRS)	45596

Full query:

```
%sql SELECT customer, SUM(PAYLOAD_MASS__KG_) AS total_payload FROM SPACEXTABLE WHERE customer="NASA (CRS)" GROUP BY customer;
```

Explanation:

SUM(PAYLOAD_MASS__KG_) AS total_payload – to extract total Payload mass

WHERE customer="NASA (CRS)" – to return only NASA records

GROUP BY customer – to return the total Payload

Total Payload mass carried by boosters launched by NASA was 45,496 kg

Average Payload Mass by F9 v1.1

```
%sql SELECT "Booster_Version", AVG(PAYLOAD_MASS__KG_) AS avg_payload FROM SPACEXTABLE WHERE "Booster_Version"="F9 v1.1"  
* sqlite:///my_data1.db  
Done.  


| Booster_Version | avg_payload |
|-----------------|-------------|
| F9 v1.1         | 2928.4      |


```

Full query:

```
%sql SELECT "Booster_Version", AVG(PAYLOAD_MASS__KG_) AS avg_payload FROM SPACEXTABLE  
WHERE "Booster_Version"="F9 v1.1" GROUP BY "Booster_Version";
```

Explanation:

AVG(PAYLOAD_MASS__KG_) AS avg_payload – to extract average Payload mass

WHERE "Booster_Version"="F9 v1.1" – to return only records with Booster Version F9 v1.1

GROUP BY "Booster_Version" – to return the average Payload

Average Payload mass carried by F9 v1.1 boosters was 2,928.4 kg

First Successful Ground Landing Date

```
%sql SELECT "Landing_Outcome", MIN("Date") AS min_date FROM SPACEXTABLE WHERE "Landing_Outcome"="Success (ground pad)" GROUP BY "Landing_Outcome";  
* sqlite:///my_data1.db  
Done.  


| Landing_Outcome      | min_date   |
|----------------------|------------|
| Success (ground pad) | 2015-12-22 |


```

Full query:

```
%sql SELECT "Landing_Outcome", MIN("Date") AS min_date FROM SPACEXTABLE WHERE "Landing_Outcome"="Success (ground pad)" GROUP BY "Landing_Outcome";
```

Explanation:

MIN("Date") AS min_date – to extract the first date

WHERE "Landing_Outcome"="Success (ground pad)" – to return only records with successful Ground Landing

GROUP BY "Landing_Outcome" – to return the date

The first successful Ground Landing was on 2015-12-22 (22 December 2015)

Successful Drone Ship Landing with Payload between 4000 and 6000

```
%sql SELECT DISTINCT "Booster_Version" FROM SPACEXTABLE 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
```

Full query:

```
%sql SELECT DISTINCT "Booster_Version" FROM SPACEXTABLE WHERE "Landing_Outcome"="Success (drone ship)" AND (PAYLOAD_MASS_KG_>4000 AND PAYLOAD_MASS_KG_<6000);
```

Explanation:

SELECT DISTINCT "Booster_Version" – to return unique Booster Versions

WHERE "Landing_Outcome"="Success (drone ship)" AND (PAYLOAD_MASS_KG_>4000 AND PAYLOAD_MASS_KG_<6000) – to return only records with successful Drone Ship landing and specified Payload Mass

Booster versions F9 FT-B1022, -B1026, -B1021.2, and -B1031.2 fulfilled the specified requirements

Total Number of Successful and Failure Mission Outcomes

```
%sql SELECT "Mission_Outcome", COUNT(*) AS count_mission_outcome FROM SPACEXTABLE GROUP BY TRIM("Mission_Outcome")
```

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

Mission_Outcome	count_mission_outcome
Failure (in flight)	1
Success	99
Success (payload status unclear)	1

Full query:

```
%sql SELECT "Mission_Outcome", COUNT(*) AS count_mission_outcome FROM SPACEXTABLE GROUP BY TRIM("Mission_Outcome");
```

Explanation:

SELECT "Mission_Outcome", COUNT() AS count_mission_outcome – to return count of Missions
GROUP BY ("Mission_Outcome"); – to group results by Mission Outcome*

In total, 1 Mission failed, 1 was successful but Payload status was unclear, and
99 Missions were successful

Boosters Carried Maximum Payload

```
%sql SELECT "Booster_Version" FROM SPACEXTABLE WHERE PAYLOAD_MASS__KG_ = (SELECT MAX(PAYLOAD_MASS__KG_) AS max_payload FROM SPACEXTABLE);  
  
* 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
```

Full query:

```
%sql SELECT "Booster_Version" FROM SPACEXTABLE WHERE PAYLOAD_MASS__KG_ =  
(SELECT MAX(PAYLOAD_MASS__KG_) AS max_payload FROM SPACEXTABLE);
```

Explanation:

SELECT "Booster_Version", – to return count of Booster Version

*WHERE PAYLOAD_MASS__KG_ = (SELECT MAX(PAYLOAD_MASS__KG_) AS max_payload
FROM SPACEXTABLE); – to return only records with maximum Payload mass*

The Booster versions listed in the table all carried the Maximum
Payload Mass

2015 Launch Records

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



| year | month | Landing_Outcome      | Booster_Version | Launch_Site |
|------|-------|----------------------|-----------------|-------------|
| 2015 | 01    | Failure (drone ship) | F9 v1.1 B1012   | CCAFS LC-40 |
| 2015 | 04    | Failure (drone ship) | F9 v1.1 B1015   | CCAFS LC-40 |


```

Full query:

```
%sql SELECT substr("Date", 1, 4) AS year, 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';
```

Explanation:

SELECT "substr("Date", 1, 4) AS year, substr("Date", 6, 2) AS month– to return year and month, as SQLite doesn't support month names

WHERE "Landing_Outcome"= "Failure (drone ship)" AND substr("Date",0,5)='2015';– to return only values where landing failed on a Drone Ship and Launch Year was 2015

Two Missions where the Landing Site was a Drone Ship failed in 2015, one in January and the other in April

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

```
%sql select count(*) co , Landing_Outcome from SPACEXTABLE where Date>= '2010-06-04' and Date <= '2017-03-20' group by Landing_Outcome order by co DESC;
```

* sqlite:///my_data1.db
Done.

co	Landing_Outcome
10	No attempt
5	Success (drone ship)
5	Failure (drone ship)
3	Success (ground pad)
3	Controlled (ocean)
2	Uncontrolled (ocean)
2	Failure (parachute)
1	Precluded (drone ship)

Full query:

%sql SELECT count() co , Landing_Outcome from SPACEXTABLE WHERE Date>= '2010-06-04' AND Date <= '2017-03-20' GROUP BY Landing_Outcome order by co DESC;*

Explanation:

SELECT count() co , Landing_Outcome – to return count as co and Landing Outcome*

WHERE Date>= '2010-06-04' AND Date <= '2017-03-20' – to return only records within set data limit

GROUP BY Landing_Outcome order by co DESC; - to group by Outcome, and rank

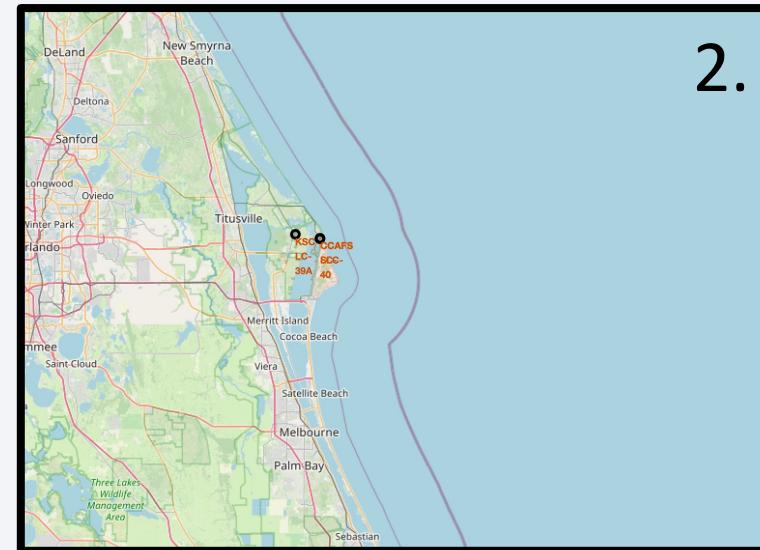
Landing Outcome ranking was as listed to the right

The background of the slide is a photograph taken from space at night. It shows the curvature of the Earth against a dark blue-black void of space. 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, the green and yellow glow of the aurora borealis is visible. The atmosphere of the Earth is thin and hazy, appearing as a light blue band near the horizon.

Section 3

Launch Sites Proximities Analysis

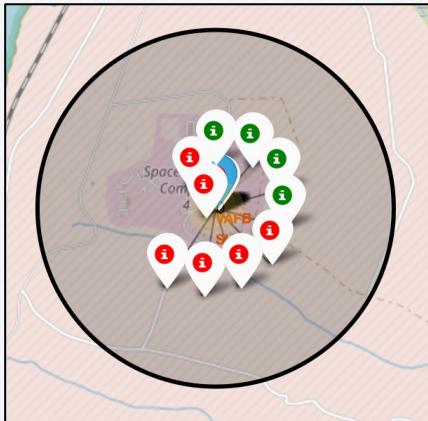
Overview of all Launch Sites



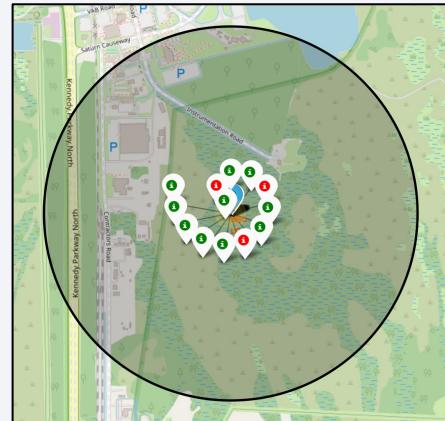
- All Launch Sites are located close to a coast line
- The Launch Sites located in Florida are closer to the equator

Launch Sites – Close-up and Mission Outcomes

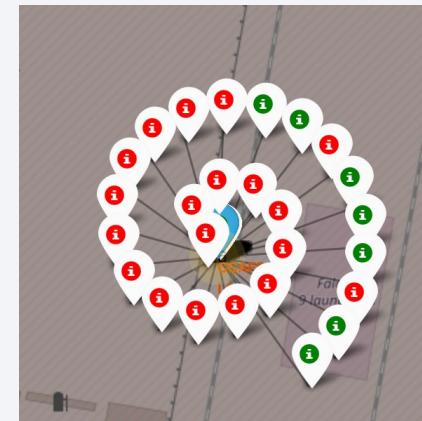
Vandenberg
Space Force Base



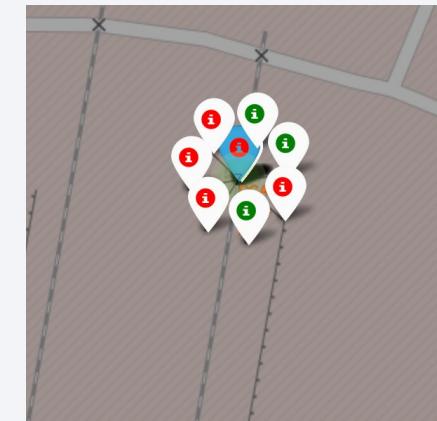
Kennedy
Space Center



Cape Canaveral
Space Launch Center

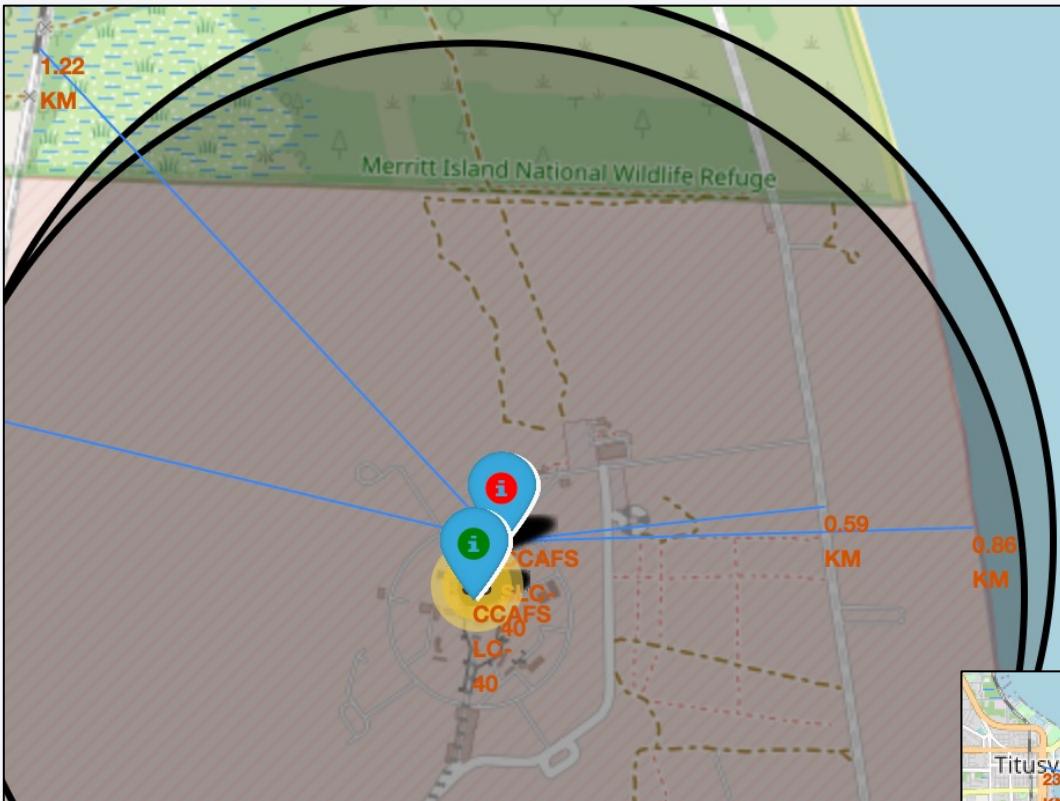


Cape Canaveral
Launch Complex



- Launch Sites are indicated by the black circle
- Successful Missions are indicated by a green marker
- Failed Missions are indicated by a red marker
- Kennedy Space Center has the highest Mission Success Rate

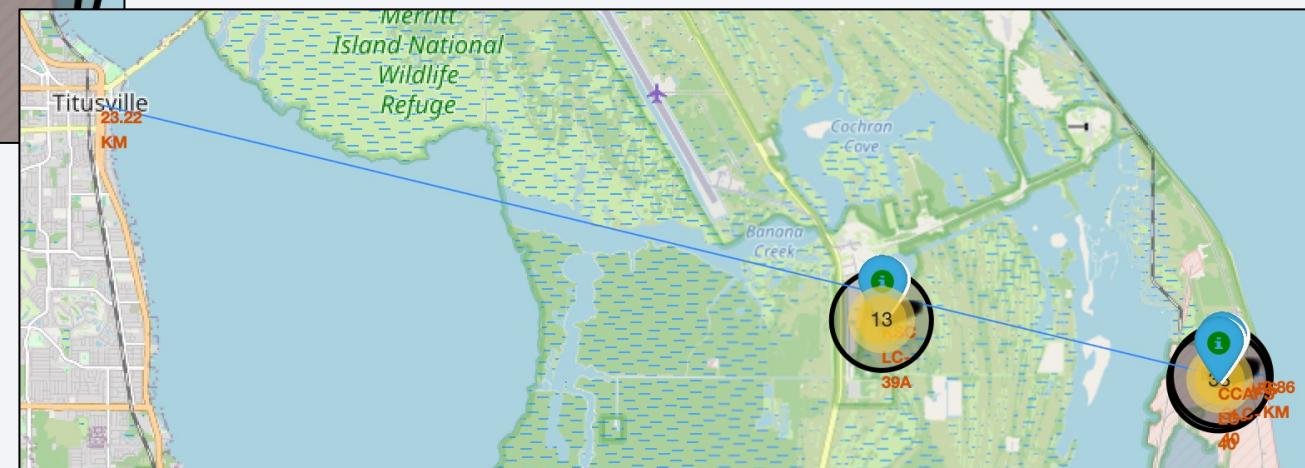
Distance between Launch Sites and Point of Interest



- The distance between the CCAFS SLC several points of interest are indicated by the blue line and distance indicator

Distances:

- Upper screenshot:
 - Highway: 0.59 km
 - Coast line: 0.86 km
 - Railway: 1.22 km
- Lower screenshot:
 - City (Titusville): 23.22 km

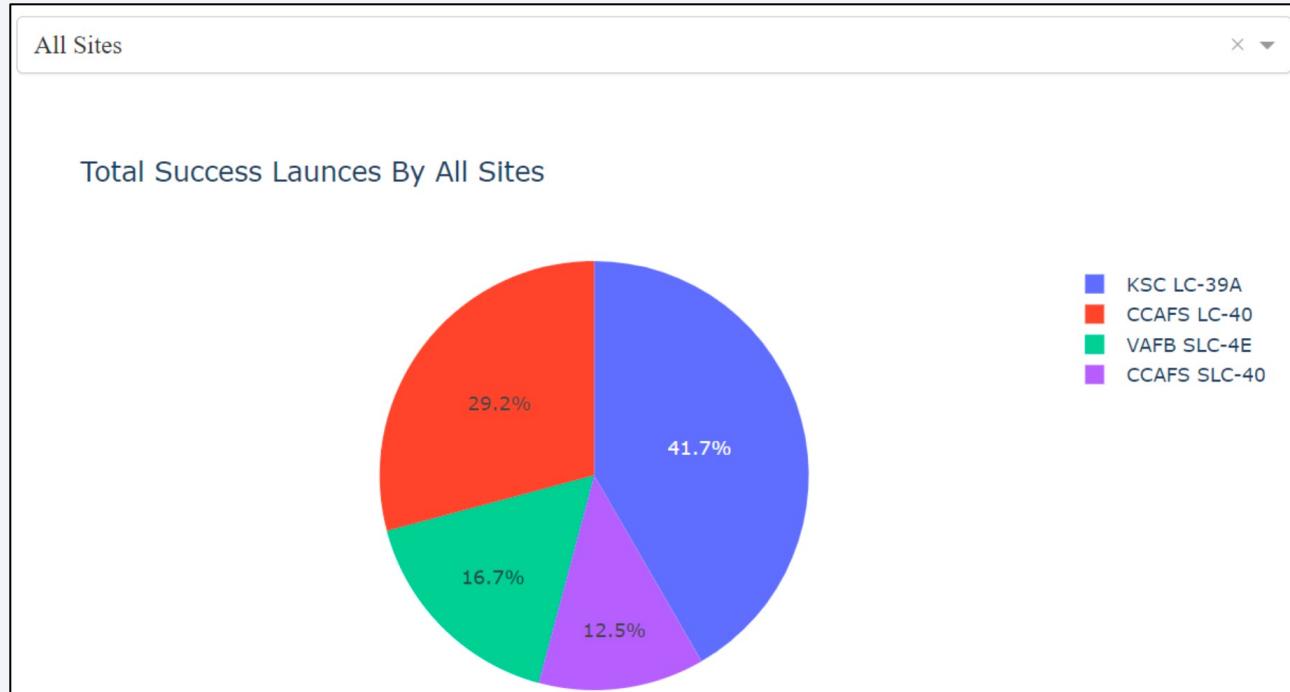


Section 4

Build a Dashboard with Plotly Dash

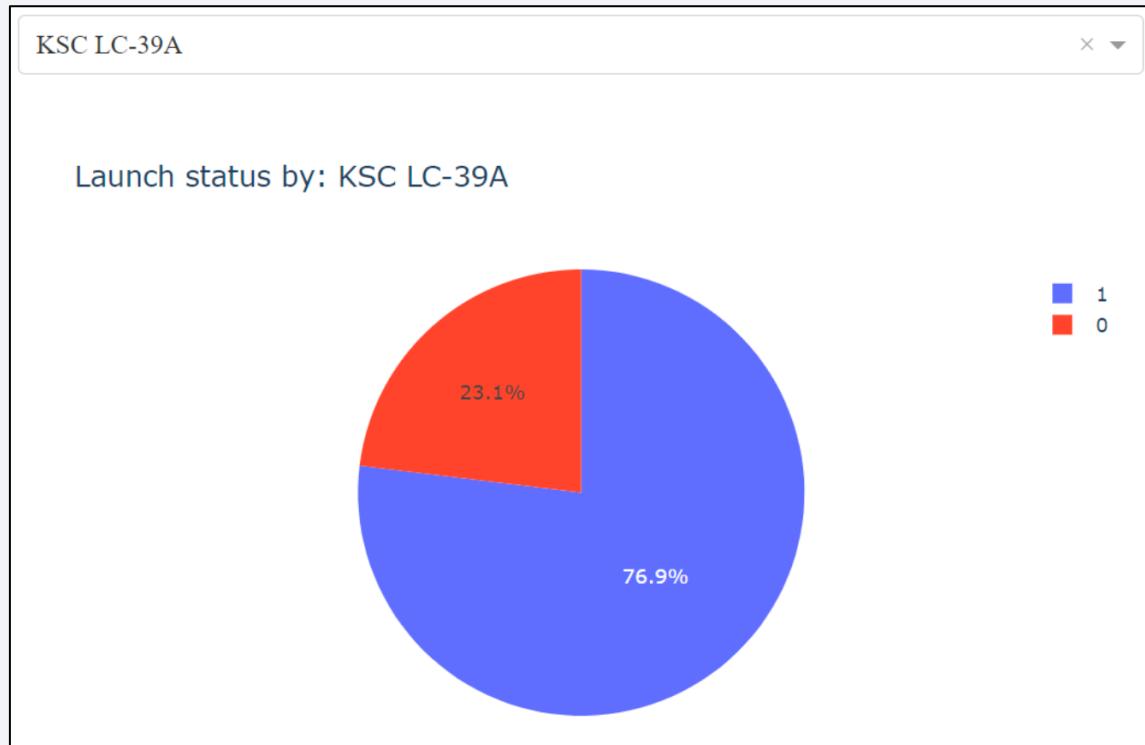


Launch Success: contribution per Launch Site



- Kennedy Space Center (KSC LC, blue) contributes most to Mission Success, with 41.7% of the successful launches being performed there
- Cape Canaveral Space Launch Center (CCAFS SLC, purple) contributed the least to Mission Success, with 12.5% of the successful launches originating from there

Launch Success: rate at the most successful Launch Site



- Kennedy Space Center (KSC LC, blue) is the most successful Launch Site
- 76.9% of the Launches was successful here (blue, Class 1)
- 23.1% of the Launches failed (red, Class 0)

Effect of Payload Mass on Launch Outcome



- Optimal Payload Mass appears to range from 2,000 to 6,000 kg

- Booster version FT appears to have the highest success rate



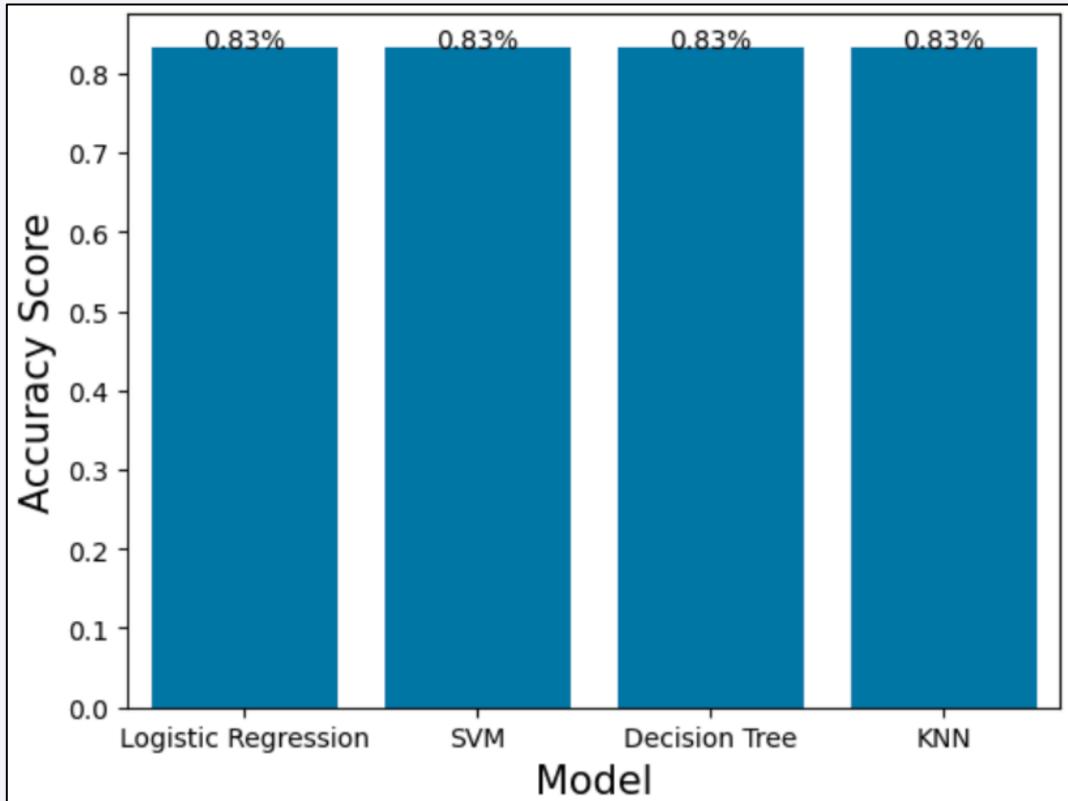
- Booster v1.1 is generally unsuccessful, even in the optimal Payload Mass range

The background of the slide features a dynamic, abstract design. It consists of several thick, curved lines that transition from a bright yellow at the top right to a deep blue at the bottom left. These lines create a sense of motion and depth, resembling a tunnel or a stylized road. The overall effect is modern and professional.

Section 5

Predictive Analysis (Classification)

Classification Accuracy

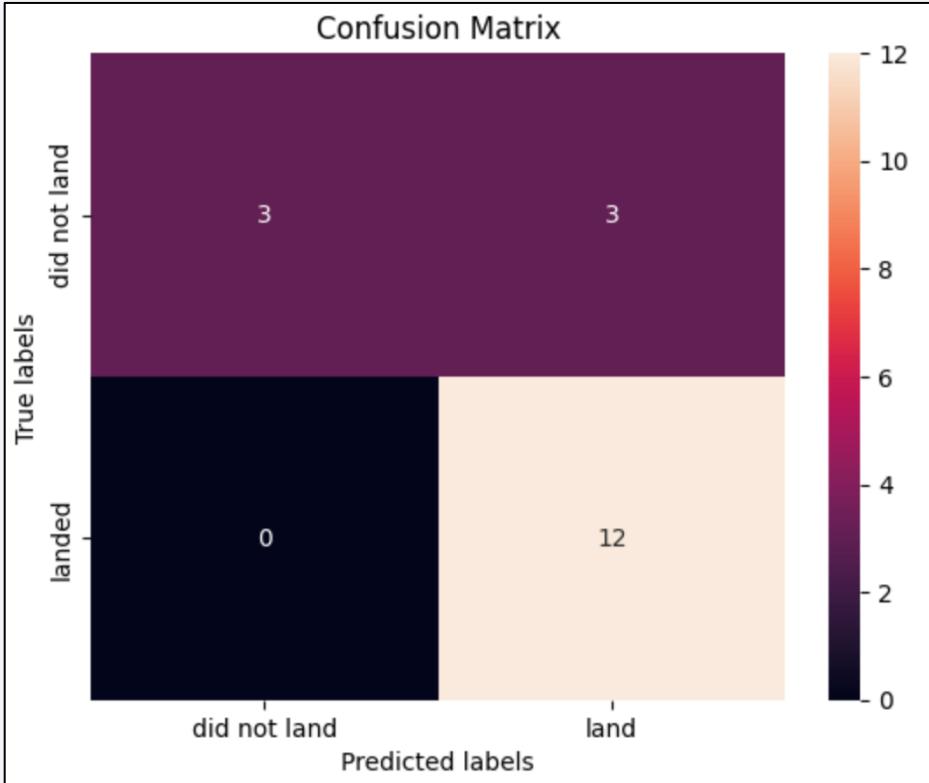


Four Predictive Models were Developed and tested:

- Logistic Regression
- Support Vector Machines
- Decision Tree
- k-Nearest Neighbours

With a Classification Accuracy of 0.833333, all predictive models performed equally well

Confusion Matrix



- All models performed equally well, and generated the same Confusion Matrix
- The Confusion Matrix of k-Nearest Neighbours is shown here
- The models work well in predicting a positive outcome, but not for the negative outcome

Conclusions

- To maintain the financial edge of SpaceX, it is vital that the Falcon 9's Reusable First Stage (RFS) lands successfully, in order to be re-used
- We here use Data Science to optimise Mission Success Rate, collecting historical data utilising the SpaceX API and web scraping, applying data wrangling, analysing data, and generating predictive models
- With time and increasing Flight Number, Mission Success Rate increased hugely, highlighting the importance of gaining experience
- The optimal Payload Mass appears to range from 2K to 6K kg, although a clear correlation could not be detected
- Kennedy Space Center is the most successful base, with a success rate of over 75%
- Prediction of Mission Failure proved challenging. Supplementing the models with additional parameters might be useful to address this issue

Thank you!

