

# Winning Space Race with Data Science

Gregory Lawson  
24/06/2023



# Outline

---

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

# Executive Summary

---

## Summary of methodologies

- Data collection/wrangling.
- Exploratory data analysis with data visualization and SQL.
- Building an interactive map with Folium.
- Building a dashboard with Dash.
- Classification (Predictive analysis).

## Summary of all results

- Exploratory data analysis results.
- Predictive analysis results.

# Introduction

---

- Project background and context

The Space Exploration Technologies Corporation (SpaceX) is an American spacecraft manufacturer, launcher, and satellite communications company headquartered in Hawthorne, California. The company was founded in 2002 by Elon Musk with the goal of reducing space transportation costs and to colonize Mars.

The company manufactures the Falcon 9, Falcon Heavy and Starship heavy-lift launch vehicles, the Cargo Dragon and Crew Dragon spacecrafts, the Starlink mega-constellation satellite and rocket engines.

The Falcon 9 rocket launches cost around 1/3<sup>rd</sup> of the competition at a cost of \$62 million per launch compared to an upward cost of \$165 million. Much of the savings is because SpaceX can reuse the first stage. Therefore, if we can determine if the first stage will land, we can determine the cost of a launch. Based on public information and machine learning models, we are going to predict if SpaceX will reuse the first stage going forward.

- Problems you want to find answers

- How do variables such as payload mass, launch site, number of flights, and orbits affect the success rate.
- Does the rate of successful landings increase over time, showing learning from trial and error.
- What is the best algorithm that can be used for binary classification in this case.

Section 1

# Methodology

# Methodology

---

## Executive Summary

- Data collection methodology:
  - Using SpaceX Rest API
  - Using Web Scrapping from Wikipedia
- Perform data wrangling
  - Filtering the data
  - Dealing with missing values
  - Using One Hot Encoding to prepare the data to a binary classification
- Perform exploratory data analysis (EDA) using visualization and SQL
- Perform interactive visual analytics using Folium and Plotly Dash
- Perform predictive analysis using classification models
  - - Building, tuning and evaluation of classification models to ensure the best results

# Data Collection

---

- Describe how data sets were collected.

Data collection process involved a combination of API requests from SpaceX REST API and Web Scraping data from a table in SpaceX's Wikipedia entry.

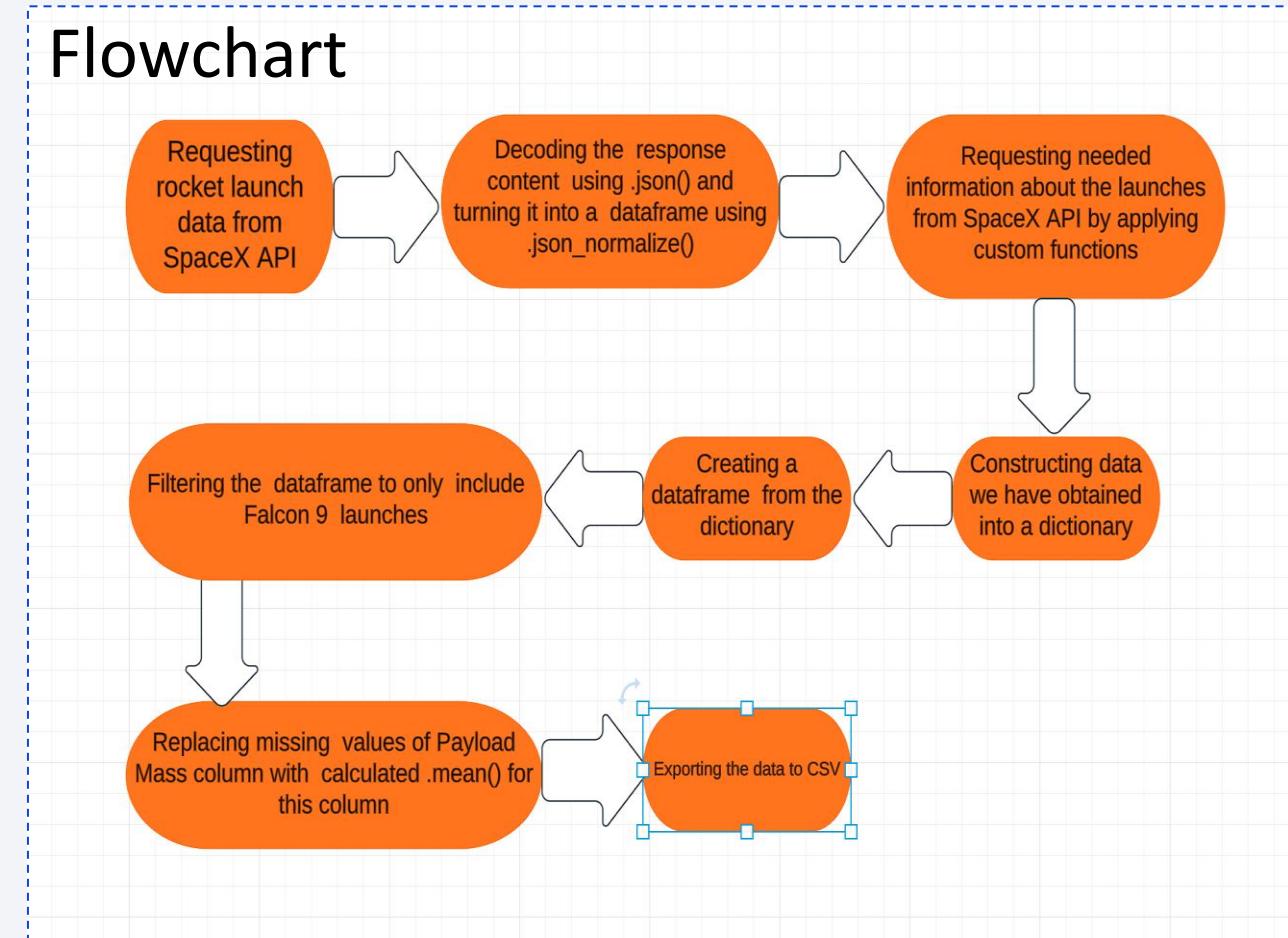
We had to use both of these data collection methods in order to get complete information about the launches for a more detailed analysis.

Data Columns are obtained by using SpaceX REST API: FlightNumber, Date, BoosterVersion, PayloadMass, Orbit, LaunchSite, Outcome, Flights, GridFins, Reused, Legs, LandingPad, Block, ReusedCount, Serial, Longitude, Latitude Data Columns are obtained by using Wikipedia Web Scraping: Flight No., Launch site, Payload, PayloadMass, Orbit, Customer, Launch outcome, Version Booster, Booster landing, Date, and Time.

# Data Collection – SpaceX API

- Present your data collection with SpaceX REST calls using key phrases and flowcharts
- Add the GitHub URL of the completed SpaceX API calls notebook (**must include completed code cell and outcome cell**), as an external reference and peer-review purpose

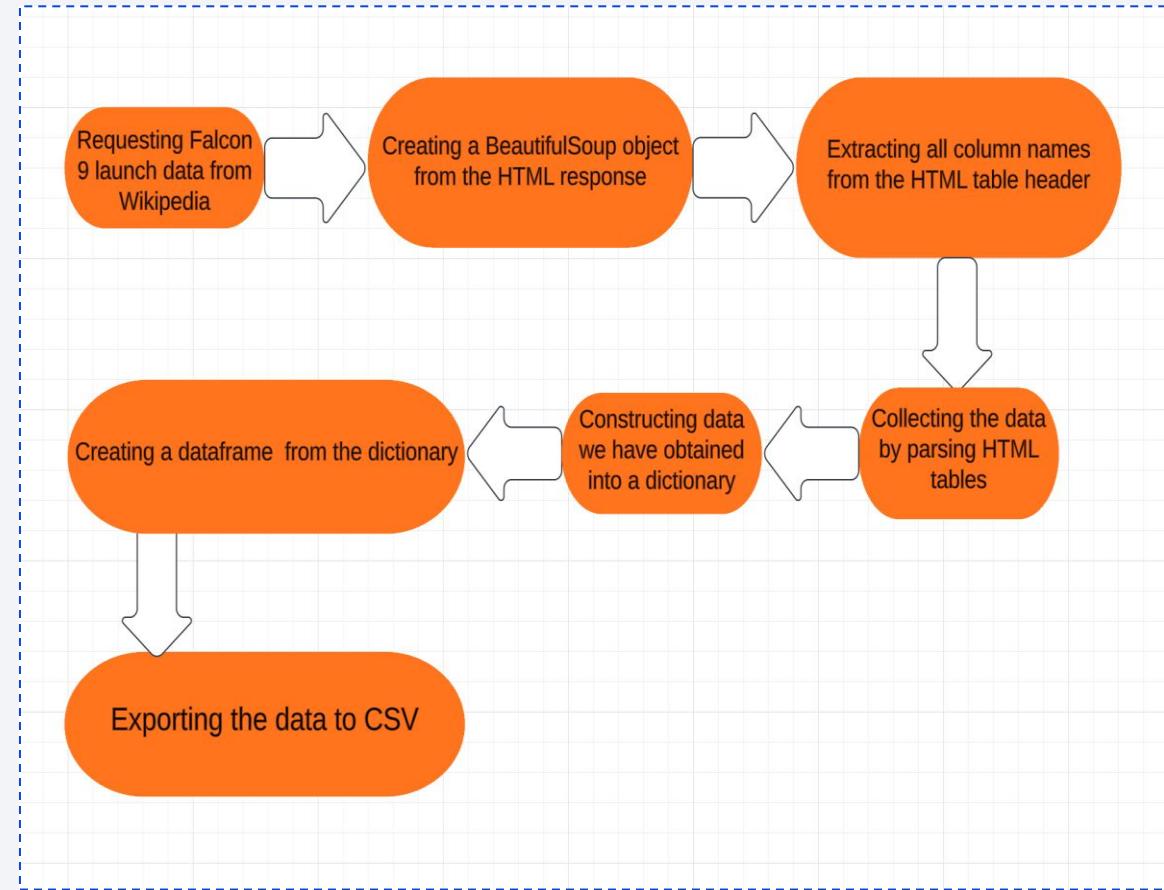
GitHub Link: [Data Collection API](#)



# Data Collection - Scraping

- Present your web scraping process using key phrases and flowcharts
- Add the GitHub URL of the completed web scraping notebook, as an external reference and peer-review purpose

GitHub Link: [Web Scraping](#)



# Data Wrangling

---

In the data set, there are several different cases where the booster did not land successfully.

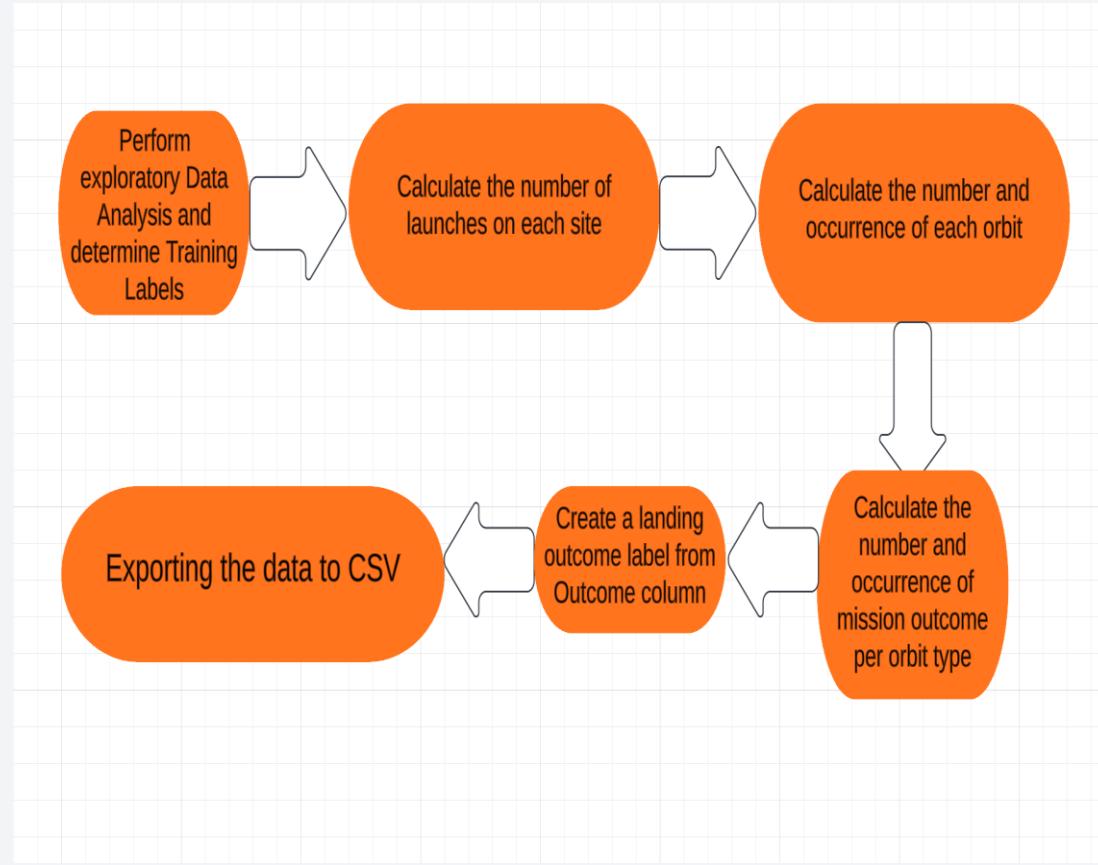
Sometimes a landing was attempted but failed due to an accident; for example, True Ocean means the mission outcome was successfully landed to a specific region of the ocean while False Ocean means the mission outcome was unsuccessfully landed to a specific region of the ocean.

True RTLS means the mission outcome was successfully landed to a ground pad  
False RTLS means the mission outcome was unsuccessfully landed to a ground pad.

True ASDS means the mission outcome was successfully landed on a drone ship  
False ASDS means the mission outcome was unsuccessfully landed on a drone ship.

We mainly convert those outcomes into Training Labels with “1” means the booster successfully landed, “0” means it was unsuccessful.

[GitHub: Data Wrangling](#)



# EDA with Data Visualization

---

Charts that were plotted:

- Flight Number vs. Payload Mass,
- Flight Number vs. Launch Site,
- Payload Mass vs. Launch Site,
- Orbit Type vs. Success Rate,
- Flight Number vs. Orbit Type,
- Payload Mass vs Orbit Type.

Success Rate Yearly Trend Scatter plots show the relationship between variables.

If a relationship exists, they could be used in machine learning model. Bar charts show comparisons among discrete categories.

The goal is to show the relationship between the specific categories being compared and a measured value. Line charts show trends in data over time (time series).

**GitHub:** [EDA with Data Viz](#)

# EDA with SQL

---

## Performed SQL queries:

- Displaying the names of the unique launch sites in the space mission
- Displaying 5 records where launch sites begin with the string ‘CCA’
- Displaying the total payload mass carried by boosters launched by NASA (CRS)
- Displaying average payload mass carried by booster version F9 v1.1
- Listing the date when the first successful landing outcome in ground pad was achieved
- Listing the names of the boosters which have success in drone ship and have payload mass greater than 4000 but less than 6000
- Listing the total number of successful and failure mission outcomes
- Listing the names of the booster versions which have carried the maximum payload mass
- Listing the failed landing outcomes in drone ship, their booster versions and launch site names for the months in year 2015 • Ranking 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.

GitHub: [EDA with SQL](#)

# Build an Interactive Map with Folium

---

## **Markers of all Launch Sites:**

- Added Marker with Circle, Popup Label and Text Label of NASA Johnson Space Center using its latitude and longitude coordinates as a start location.
- Added Markers with Circle, Popup Label and Text Label of all Launch Sites using their latitude and longitude coordinates to show their geographical locations and proximity to Equator and coasts.

## **Coloured Markers of the launch outcomes for each Launch Site:**

- Added coloured Markers of success (Green) and failed (Red) launches using Marker Cluster to identify which launch sites have relatively high success rates.

## **Distances between a Launch Site to its proximities:**

- Added coloured Lines to show distances between the Launch Site KSC LC-39A (as an example) and its proximities like Railway, Highway, Coastline and Closest City.

**GitHub:** [Interactive Map with Folium](#)

# Build a Dashboard with Plotly Dash

---

## Launch Sites Dropdown List:

- Added a dropdown list to enable Launch Site selection.

## Pie Chart showing Success Launches (All Sites/Certain Site):

- Added a pie chart to show the total successful launches count for all sites and the Success vs. Failed counts for the site, if a specific Launch Site was selected.

## Slider of Payload Mass Range:

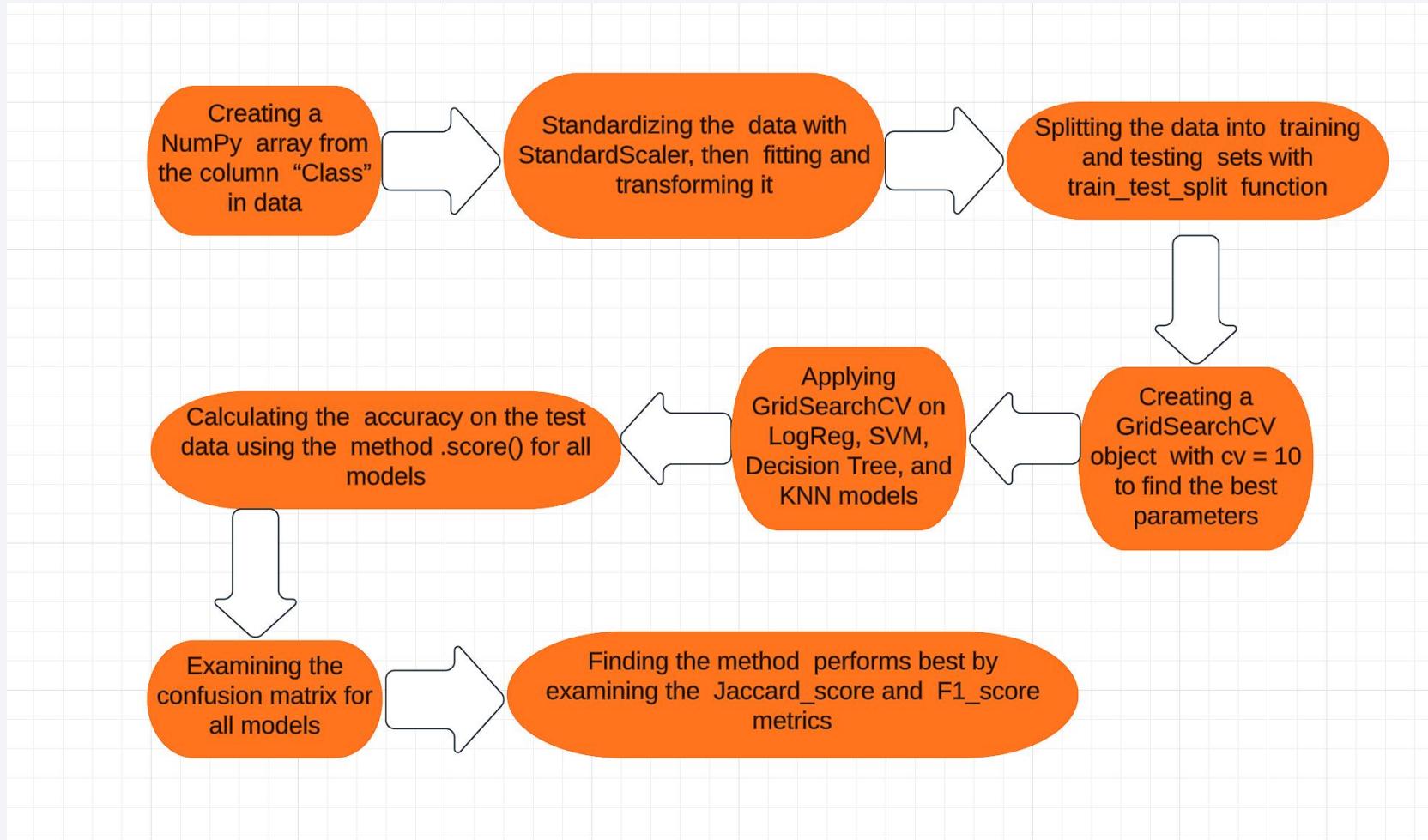
- Added a slider to select Payload range.

## Scatter Chart of Payload Mass vs. Success Rate for the different Booster Versions:

- Added a scatter chart to show the correlation between Payload and Launch Success.

[GitHub: Dashboard](#)

# Predictive Analysis (Classification)



GitHub: [Machine Learning Prediction](#)

# Results

---

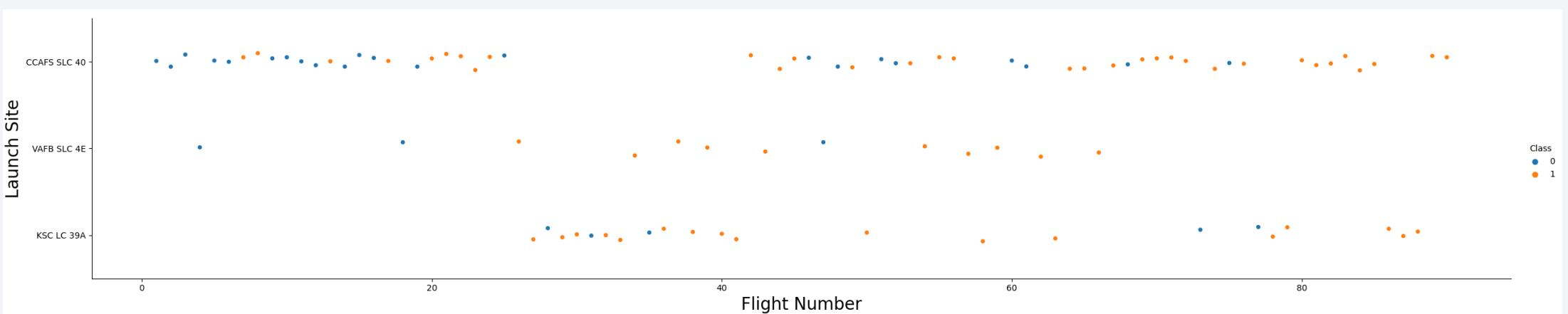
- Exploratory data analysis results
- Interactive analytics demo in screenshots
- Predictive analysis results

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

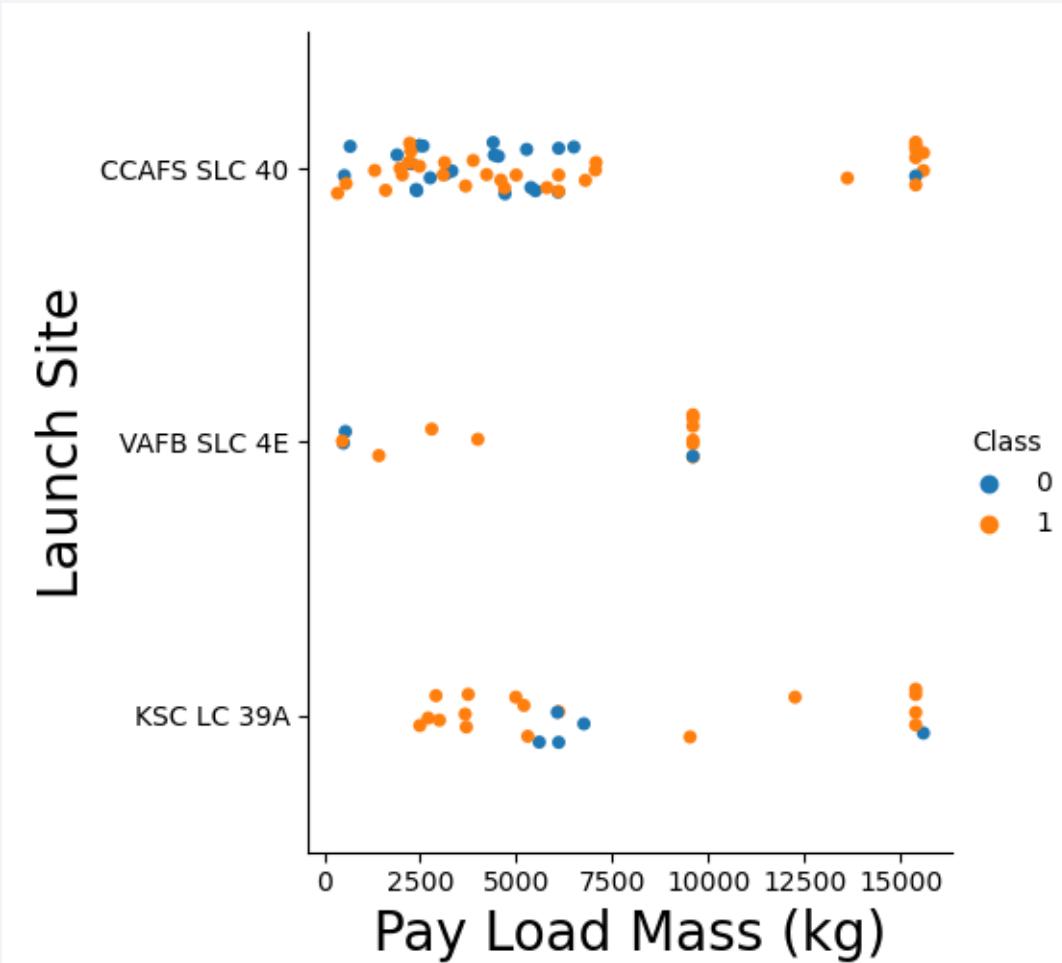
# Flight Number vs. Launch Site



## Explanation:

- The earliest flights all failed while the latest flights all succeeded.
- The CCAFS SLC 40 launch site has about a half of all launches.
- VAFB SLC 4E and KSC LC 39A have higher success rates.
- It can be assumed that each new launch has a higher rate of success.

# Payload vs. Launch Site



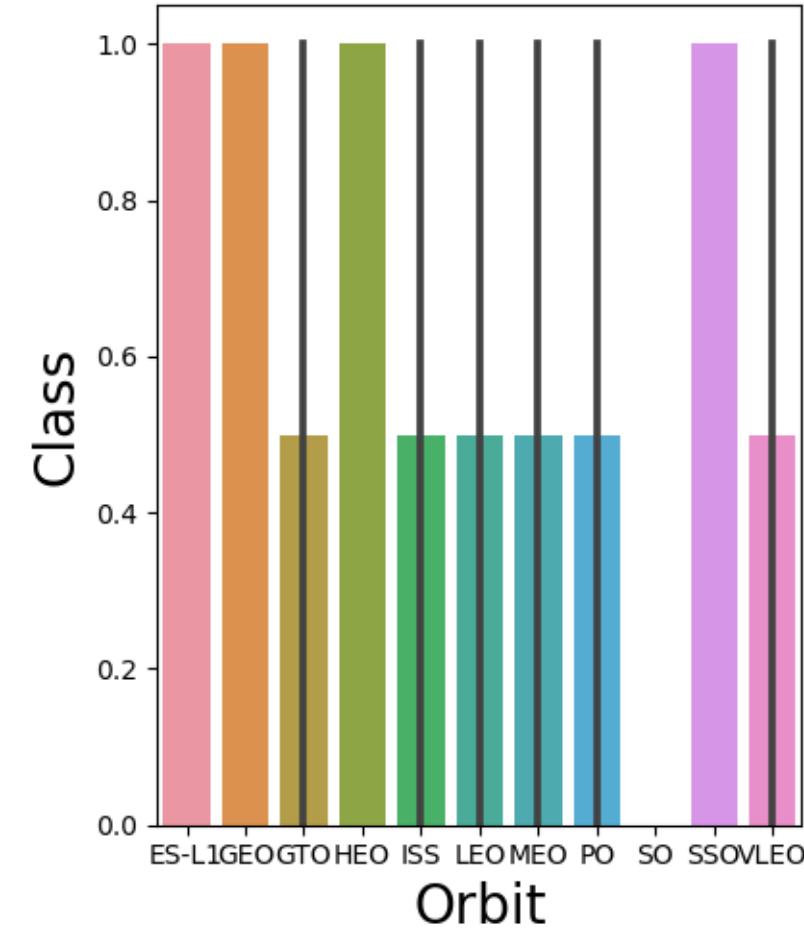
## Explanation:

- For every launch site the higher the payload mass, the higher the success rate.
- Most of the launches with payload mass over 7000 kg were successful.
- KSC LC 39A has a 100% success rate for payload mass under 5500 kg too

# Success Rate vs. Orbit Type

## Explanation:

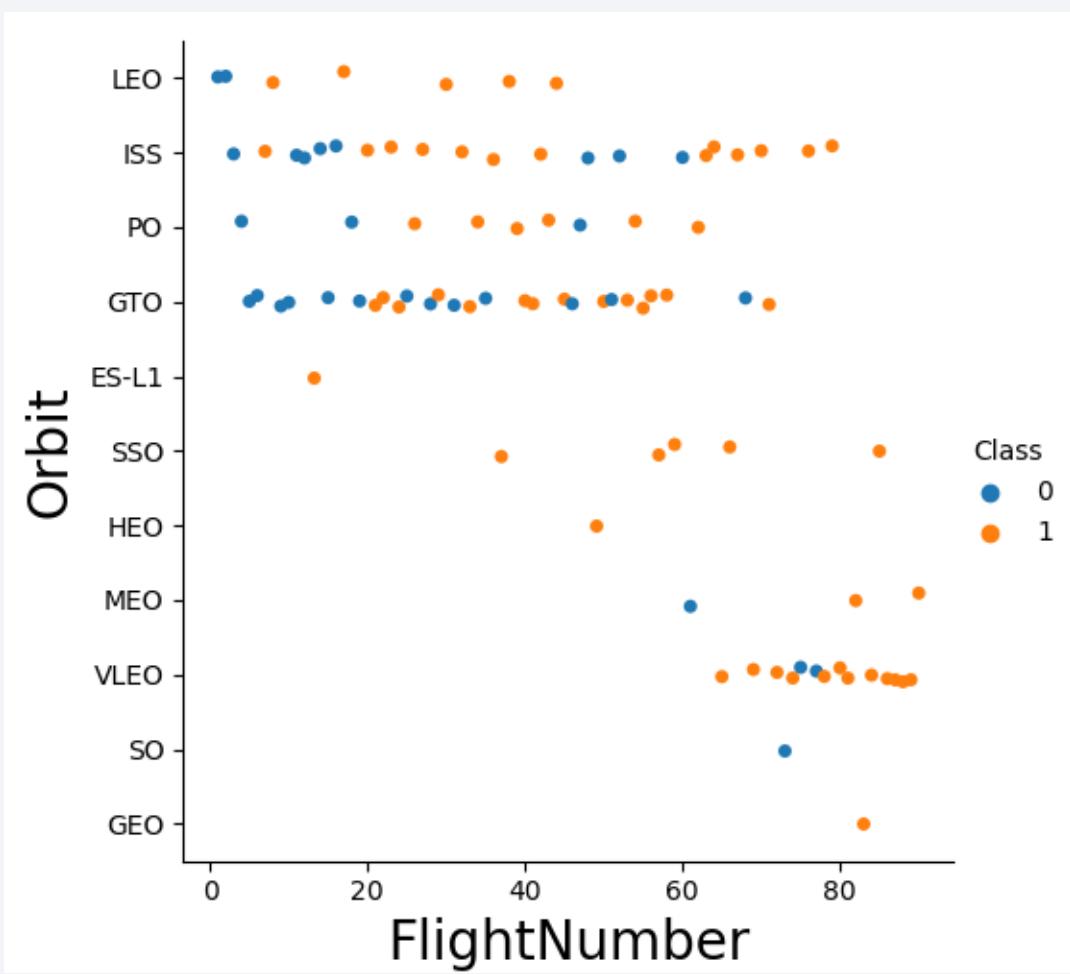
- Orbit types with 100% success rate: - ES-L1, GEO, HEO, SSO
- Orbit types with 0% success rate: - SO
- Orbit types with success rate between 50% and 85%: - GTO, ISS, LEO, MEO, PO



# Flight Number vs. Orbit Type

## Explanation:

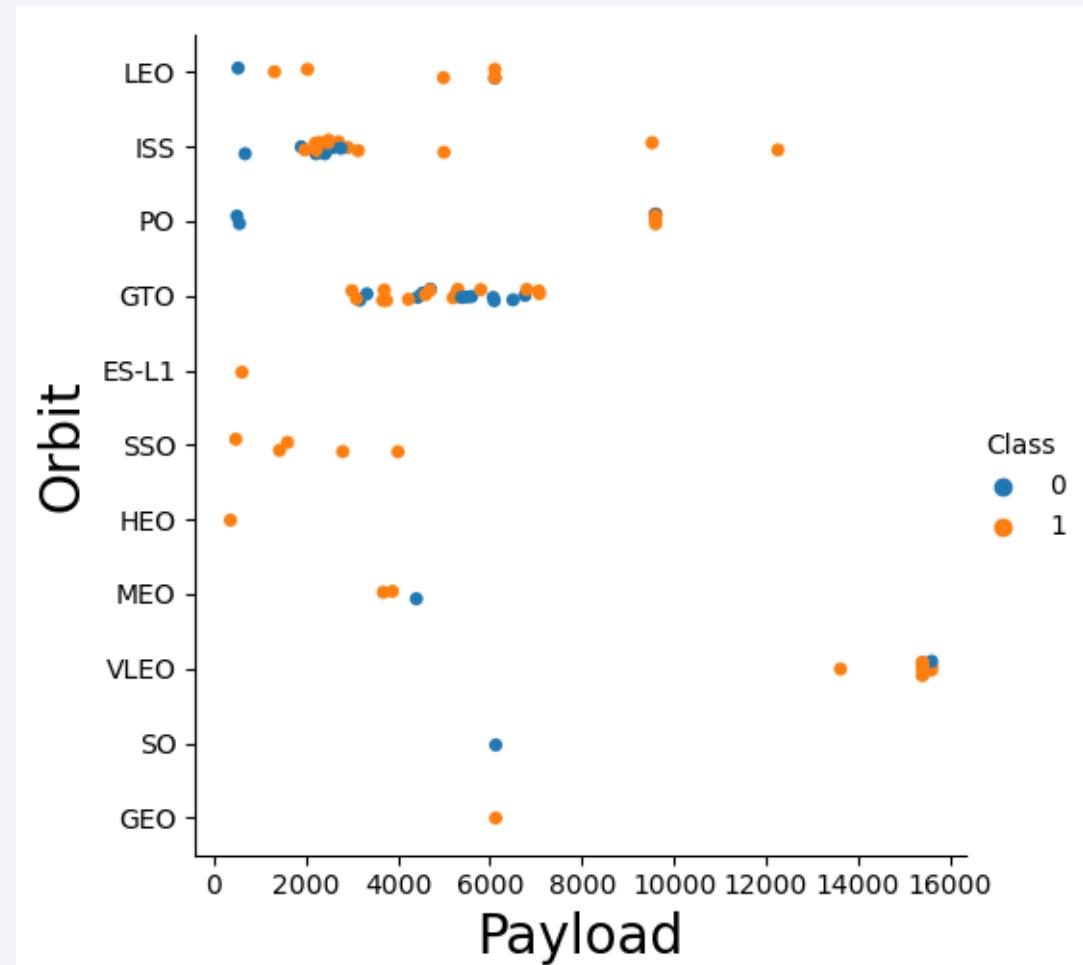
- In the LEO orbit the Success appears related to the number of flights; on the other hand, there seems to be no relationship between flight number when in GTO orbit.



# Payload vs. Orbit Type

## Explanation:

- Heavy payloads have a negative influence on GTO orbits and positive on GTO and Polar LEO (ISS) orbits.

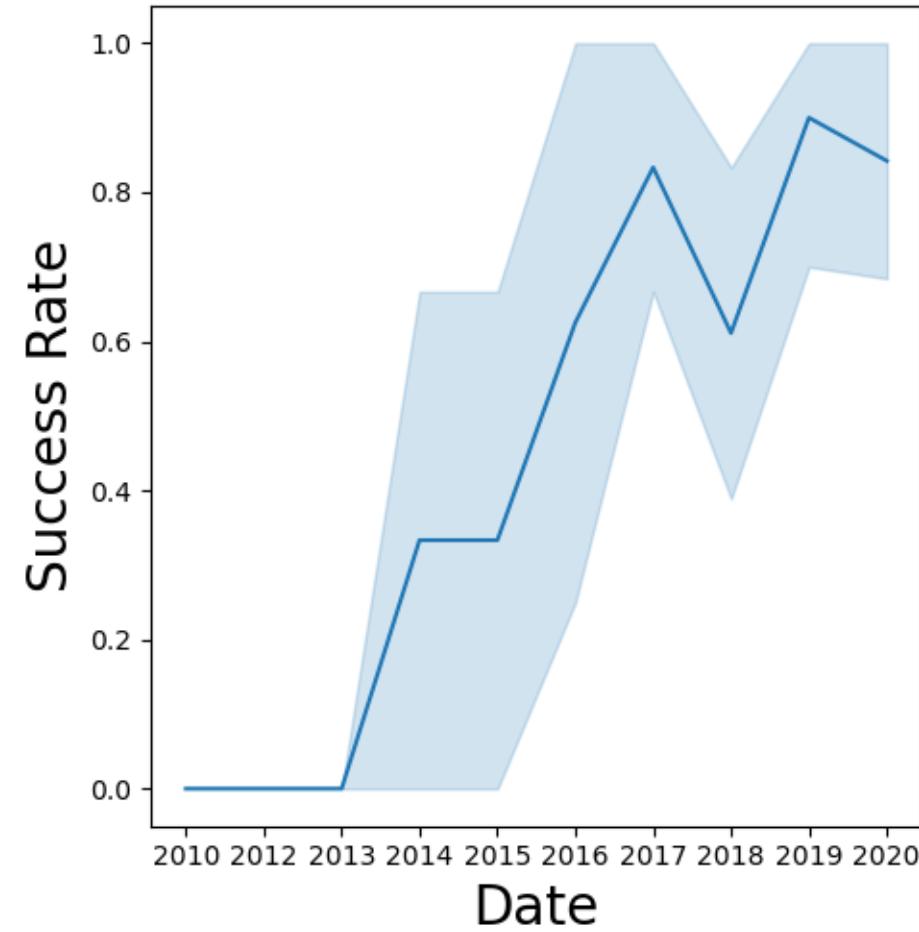


# Launch Success Yearly Trend

---

## Explanation:

- The success rate since 2013 kept increasing till 2017, with a dip in 2018 before continuing it's success.



# All Launch Site Names

---

```
In [36]: %%sql  
SELECT DISTINCT("Launch_Site") FROM SPACEXTBL;
```

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

```
Out[36]: Launch_Site
```

```
CCAFS LC-40
```

```
VAFB SLC-4E
```

```
KSC LC-39A
```

```
CCAFS SLC-40
```

## Explanation:

- Displaying the names of the unique launch sites in the space mission.

# Launch Site Names Begin with 'CCA'

```
In [37]: %%sql
SELECT * FROM SPACEXTBL WHERE LAUNCH_SITE LIKE 'CCA%' LIMIT 5;

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

Date	Time (UTC)	Booster_Version	Launch_Site	Payload	PAYLOAD_MASS_KG_	Orbit	Customer	Mission_Outcome	Landing_Outcome
06/04/2010	18:45:00	F9 v1.0 B0003	CCAFS LC-40	Dragon Spacecraft Qualification Unit	0.0	LEO	SpaceX	Success	Failure (parachute)
12/08/2010	15:43:00	F9 v1.0 B0004	CCAFS LC-40	Dragon demo flight C1, two CubeSats, barrel of Brouere cheese	0.0	LEO (ISS)	NASA (COTS) NRO	Success	Failure (parachute)
22/05/2012	7:44:00	F9 v1.0 B0005	CCAFS LC-40	Dragon demo flight C2	525.0	LEO (ISS)	NASA (COTS)	Success	No attempt
10/08/2012	0:35:00	F9 v1.0 B0006	CCAFS LC-40	SpaceX CRS-1	500.0	LEO (ISS)	NASA (CRS)	Success	No attempt
03/01/2013	15:10:00	F9 v1.0 B0007	CCAFS LC-40	SpaceX CRS-2	677.0	LEO (ISS)	NASA (CRS)	Success	No attempt

## Explanation:

- Displaying 5 records where launch sites begin with the string 'CCA'.

# Total Payload Mass

---

## Explanation:

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

```
In [40]: %%sql
SELECT sum (PAYLOAD_MASS__KG_) AS payloadmasskg FROM SPACEXTBL WHERE Customer = "NASA (CRS)"
```

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

```
Out[40]: payloadmasskg
```

```
45596.0
```

# Average Payload Mass by F9 v1.1

---

In [41]:

```
%%sql
SELECT avg(PAYLOAD_MASS__KG_) AS payloadmasskg FROM SPACEXTBL WHERE booster_version = 'F9 v1.1';

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

Out[41]: [payloadmasskg](#)

2928.4

## Explanation:

- Displaying average payload mass carried by booster version F9 v1.1.

# First Successful Ground Landing Date

---

## Explanation:

- Listing the date when the first successful landing outcome in ground pad was achieved.

In [44]:

```
%%sql
SELECT min (DATE) FROM SPACEXTBL WHERE Landing_Outcome = "Success (ground pad)"
```

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

Out[44]: min (DATE)

01/08/2018

# Successful Drone Ship Landing with Payload between 4000 and 6000

---

## Explanation:

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

```
In [45]: %%sql
SELECT BOOSTER_VERSION FROM SPACEXTBL WHERE Landing_Outcome = 'Success (drone ship)' AND PAYLOAD_MASS_KG_ > 4000 AND PAYLOAD_MASS_KG_ < 6000
* sqlite:///my_data1.db
Done.

Out[45]: Booster_Version
F9 FT B1022
F9 FT B1026
F9 FT B1021.2
F9 FT B1031.2
```

# Total Number of Successful and Failure Mission Outcomes

---

## Explanation:

- Listing the total number of successful and failure mission outcomes.

```
In [46]: %%sql  
SELECT COUNT(MISSION_OUTCOME) FROM SPACEXTBL WHERE MISSION_OUTCOME = 'Success' or MISSION_OUTCOME = 'Failure (in flight);
```

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

```
Out[46]: COUNT(MISSION_OUTCOME)
```

---

```
99
```

# Boosters Carried Maximum Payload

---

## Explanation:

- Listing the names of the booster versions which have carried the maximum payload mass.

```
In [47]: %%sql
SELECT BOOSTER_VERSION FROM SPACEXTBL WHERE PAYLOAD_MASS_KG_ = (SELECT max(PAYLOAD_MASS_KG_) FROM SPACEXTBL);

* sqlite:///my_data1.db
Done.

Out[47]: 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

---

## Explanation:

- Listing the failed landing outcomes in drone ship, their booster versions and launch site names for the months in year 2015.

```
In [50]: %%sql
SELECT Landing_Outcome, BOOSTER_VERSION, LAUNCH_SITE FROM SPACEXTBL WHERE Landing_Outcome = 'Failure (drone ship)' AND sub
* sqlite:///my_data1.db
Done.

Out[50]: 

| 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

---

## Explanation:

- Ranking 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.

```
In [51]: %%sql
SELECT COUNT(Landing_Outcome) FROM SPACEXTBL WHERE Landing_Outcome='Failure (drone ship)' OR Landing_Outcome='Success (gro
* sqlite:///my_data1.db
Done.

Out[51]: COUNT(Landing_Outcome)
0
```

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. Numerous glowing yellow and white points represent city lights, concentrated in coastal and urban areas. In the upper right quadrant, there are bright green and yellow bands of light, likely the Aurora Borealis or Australis. The overall atmosphere is dark and mysterious.

Section 3

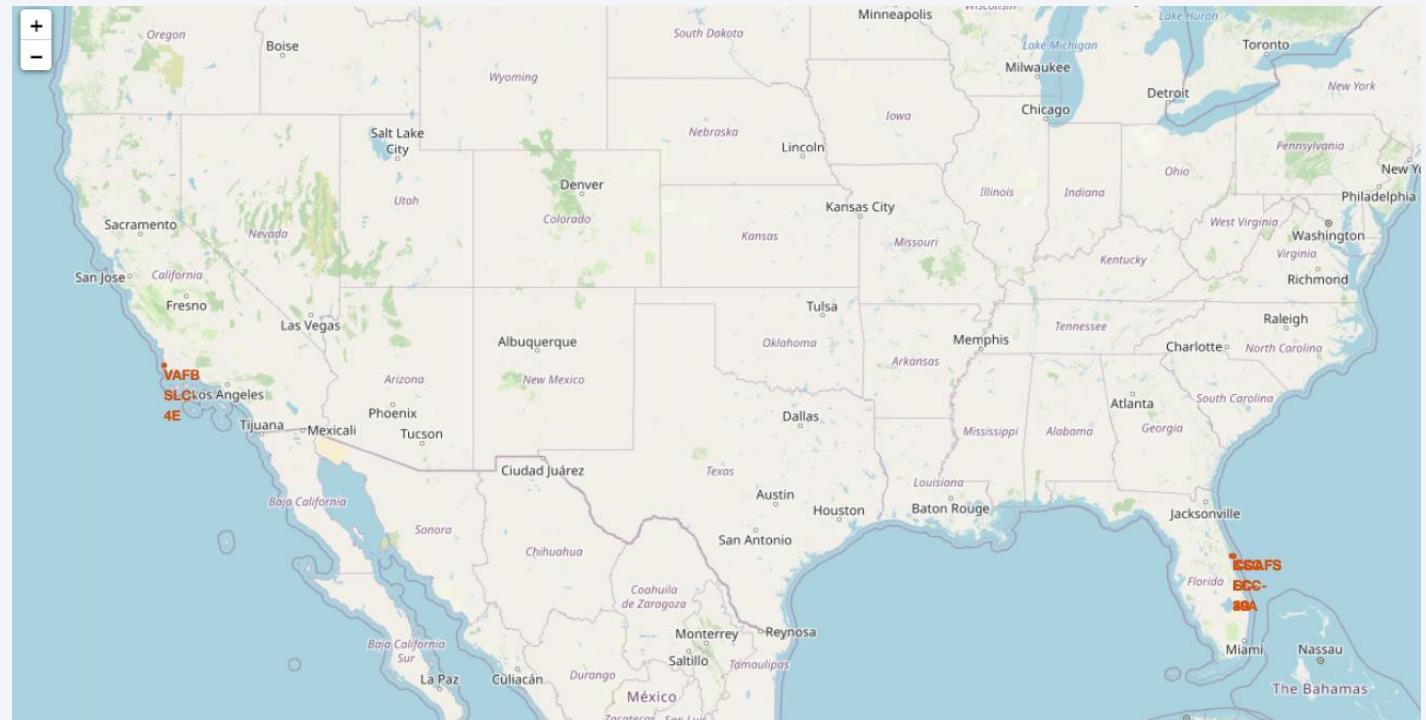
# Launch Sites Proximities Analysis

# All launch sites' location markers on a global map

---

## Explanation:

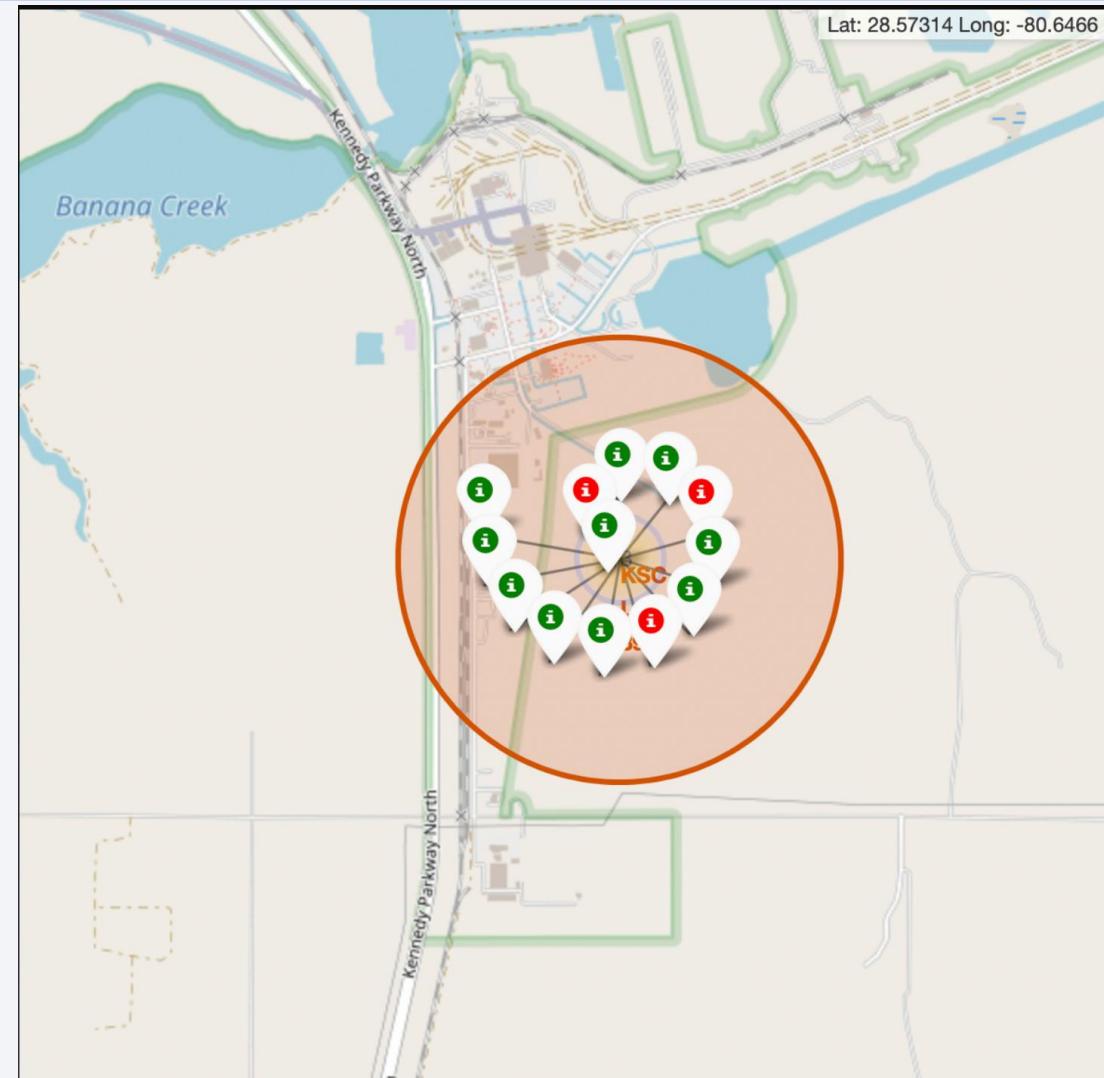
- Most of Launch sites are in proximity to the Equator line. The land is moving faster at the equator than any other place on the surface of the Earth. Anything on the surface of the Earth at the equator is already moving at 1670 km/hour. If a ship is launched from the equator it goes up into space, and it is also moving around the Earth at the same speed it was moving before launching. This is because of inertia. This speed will help the spacecraft keep up a good enough speed to stay in orbit.
- All launch sites are in very close proximity to the coast, while launching rockets towards the ocean it minimises the risk of having any debris dropping or exploding near people.



# Launch outcomes

## Explanation:

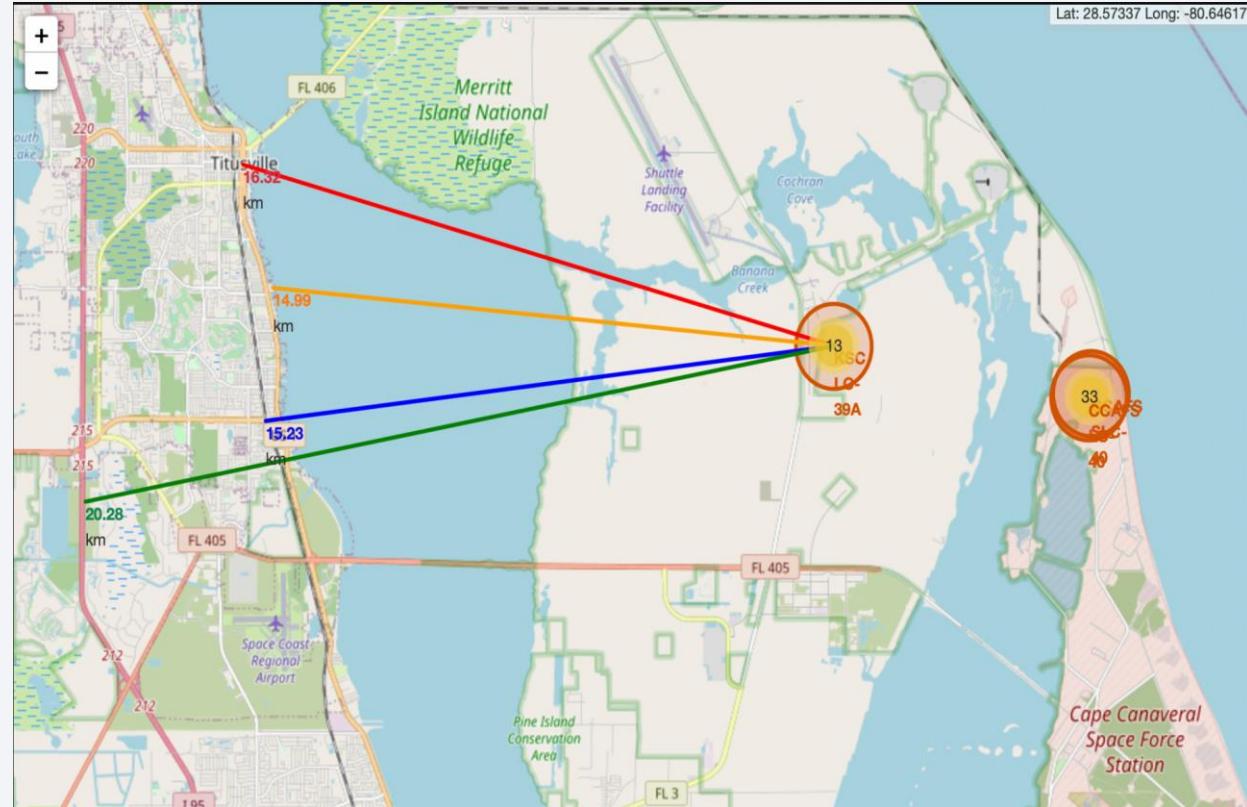
- From the colour-labeled markers we should be able to easily identify which launch sites have relatively high success rates.
  - Green Marker = Successful Launch
  - Red Marker = Failed Launch
- Launch Site KSC LC-39A has a very high Success Rate.

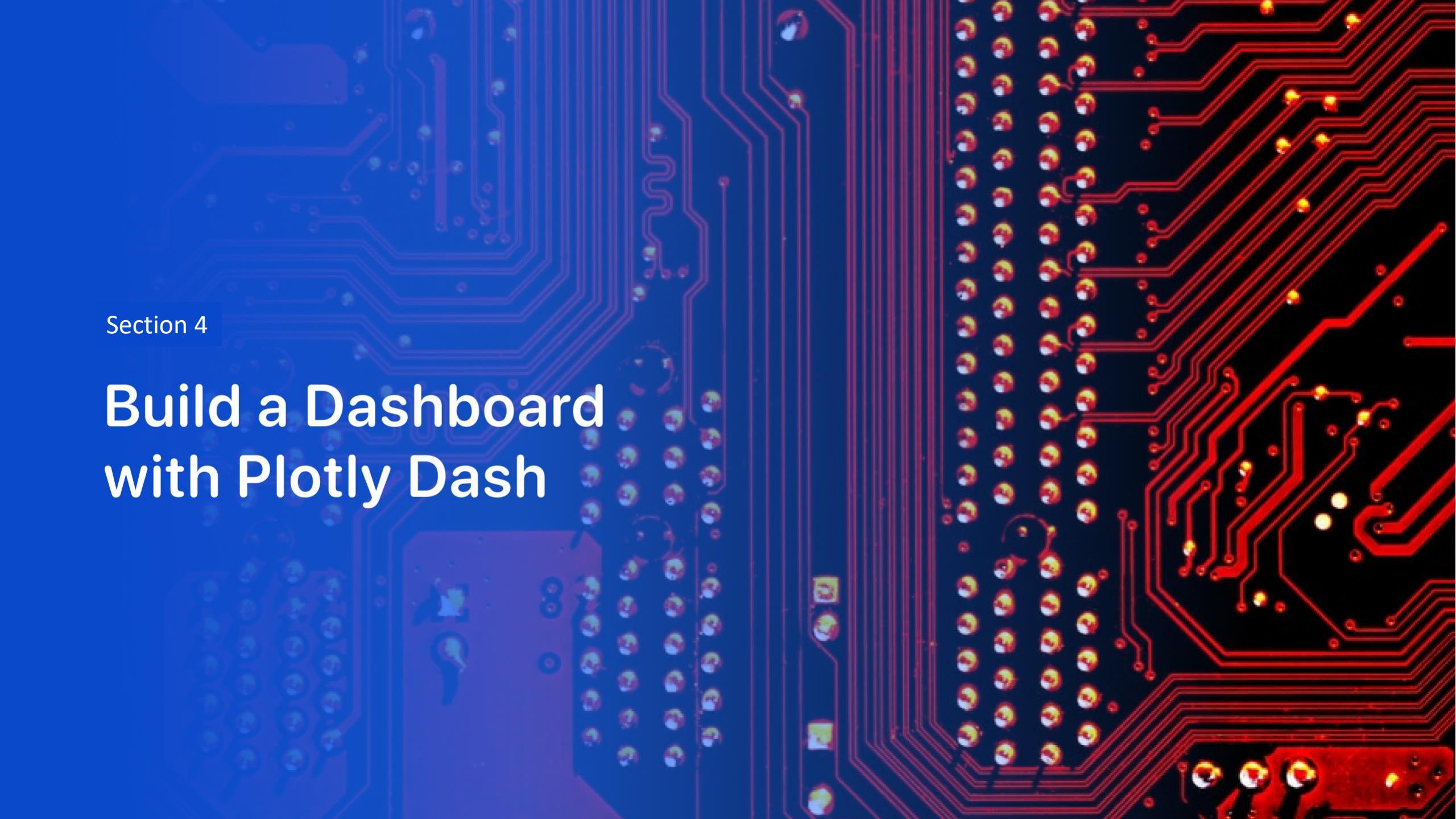


# launch site KSC LC-39A

## Explanation:

- From the visual analysis of the launch site KSC LC-39A we can clearly see that it is:
  - relative close to railway (15.23 km)
  - relative close to highway (20.28 km)
  - relative close to coastline (14.99 km)
- Also the launch site KSC LC-39A is relative close to its closest city Titusville (16.32 km).
- Failed rocket with its high speed can cover distances like 15-20 km in few seconds. It could be potentially dangerous to populated areas.





Section 4

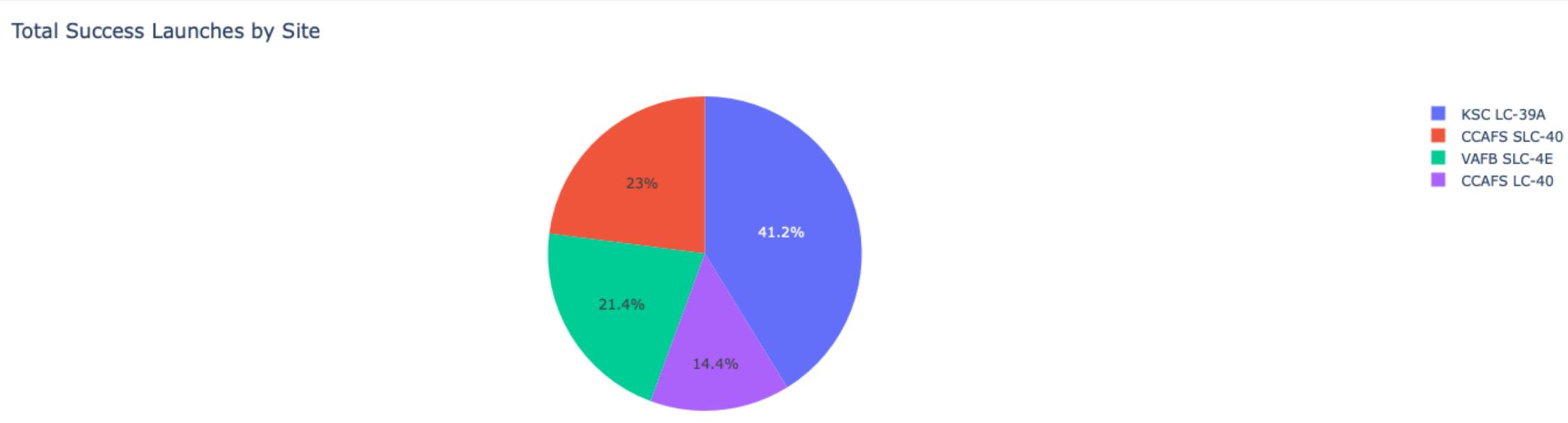
# Build a Dashboard with Plotly Dash

# Launch success count for all sites

---

## Explanation:

- The chart clearly shows that from all the sites, KSC LC-39A has the most successful launches.



# Highest launch success ratio

---

## Explanation:

- KSC LC-39A has the highest launch success rate (76.9%) with 10 successful and only 3 failed landings.

Total Success Launches for Site KSC LC-39A



# Payload vs. Launch Outcome

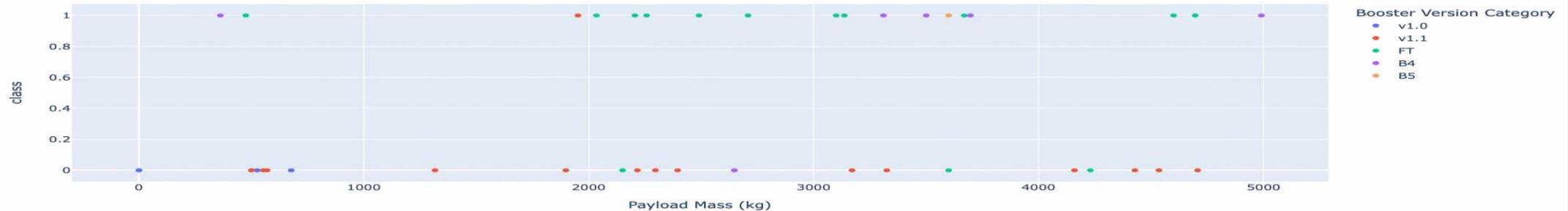
## Explanation:

- The charts show that payloads between 2000 and 5500 kg have the highest success rate.

Payload range (Kg):



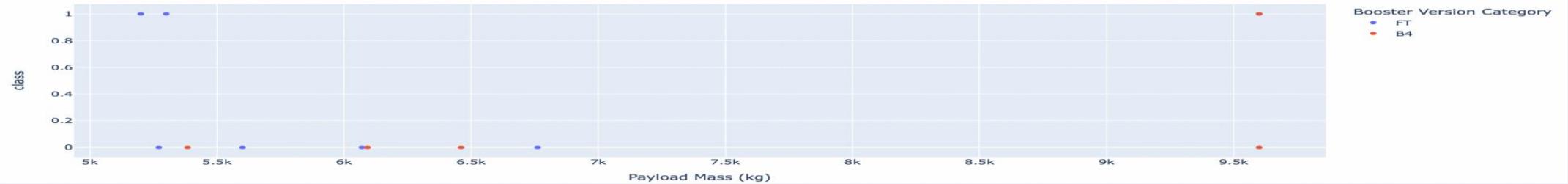
Correlation Between Payload and Success for All Sites



Payload range (Kg):



Correlation Between Payload and Success for All Sites



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 landscape. The overall effect is modern and professional.

Section 5

# Predictive Analysis (Classification)

# Classification Accuracy

---

## Explanation:

- Based on the scores of the Test Set, we can not confirm which method performs best.
- Same Test Set scores may be due to the small test sample size (18 samples). Therefore, we tested all methods based on the whole Dataset.
- The scores of the whole Dataset confirm that the best model is the Decision Tree Model. This model has not only higher scores, but also the highest accuracy.

Test set accuracy & scores

	LogReg	SVM	Tree	KNN
Jaccard_Score	0.800000	0.800000	0.800000	0.800000
F1_Score	0.888889	0.888889	0.888889	0.888889
Accuracy	0.833333	0.833333	0.833333	0.833333

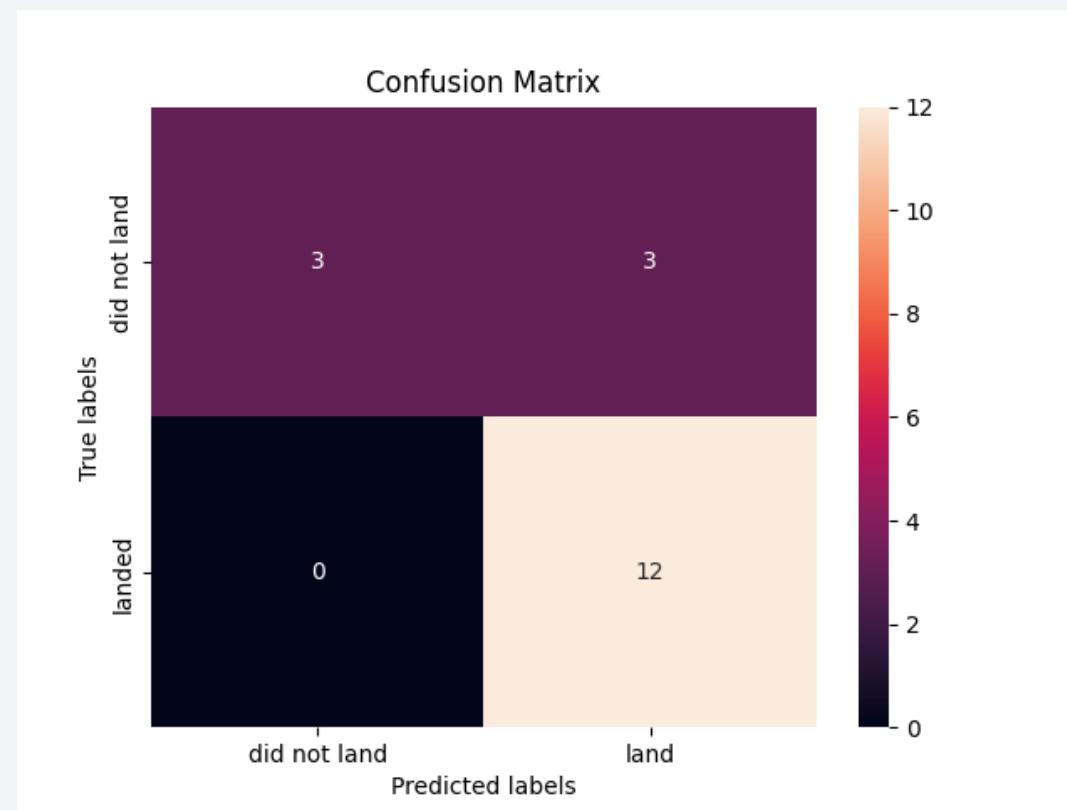
Entire data set accuracy & scores

	LogReg	SVM	Tree	KNN
Jaccard_Score	0.833333	0.845070	0.882353	0.819444
F1_Score	0.909091	0.916031	0.937500	0.900763
Accuracy	0.866667	0.877778	0.911111	0.855556

# Confusion Matrix

## Explanation:

- Examining the confusion matrix, we see that logistic regression can distinguish between the different classes. We see that the major problem is false positives.



# Conclusions

---

- The Decision Tree Model is the best algorithm for this dataset.
- Launches with a low payload mass show better results than launches with a larger payload mass.
- Most of launch sites are in proximity to the Equator line and all the sites are in very close proximity to the coast.
- The success rate of launches increases over the years, showing that the company is consistently learning from their experience.
- KSC LC-39A has the highest success rate of the launches from all the sites.
- Orbits ES-L1, GEO, HEO and SSO have 100% success rate.

Thank you!

