

# Winning Space Race with Data Science

Sanjeev Sreedhar

31/07/2025



# Outline

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



# Executive Summary

## ➤ Summary of Methodologies

- Data Collection
- Data Wrangling
- Exploratory Data Analysis with Data Visualization
- Exploratory Data Analysis with SQL
- Building an interactive map with Folium
- Building a Dashboard with Plotly Dash
- Predictive analysis(Classification)

## ➤ Summary of all results

- Exploratory Data Analysis results
- Interactive analytics demo
- Predictive analysis results



# Introduction

## Project background and context:

**SpaceX** is the most successful company of the commercial space age, making space travel affordable. The company advertises **Falcon 9** rocket launches on its website, with a cost of **62 million dollars**; other providers cost upward of **165 million dollars** each, 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.

## Problems you want to find answers:

- I. How do variables such as payload mass, launch site, number of flights and orbits affect the success of the first stage landing?
- II. Does the rate of successful landings increase over the years?
- III. 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**



# Data Collection

- 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, Time.



# Data Collection – SpaceX API

Requesting rocket launch data from SpaceX API

Decoding the response content using `.json()` and turning it into a dataframe using `.json_normalize()`

Requesting needed information about the launches from SpaceX API by applying custom functions

Constructing data we have obtained into a dictionary

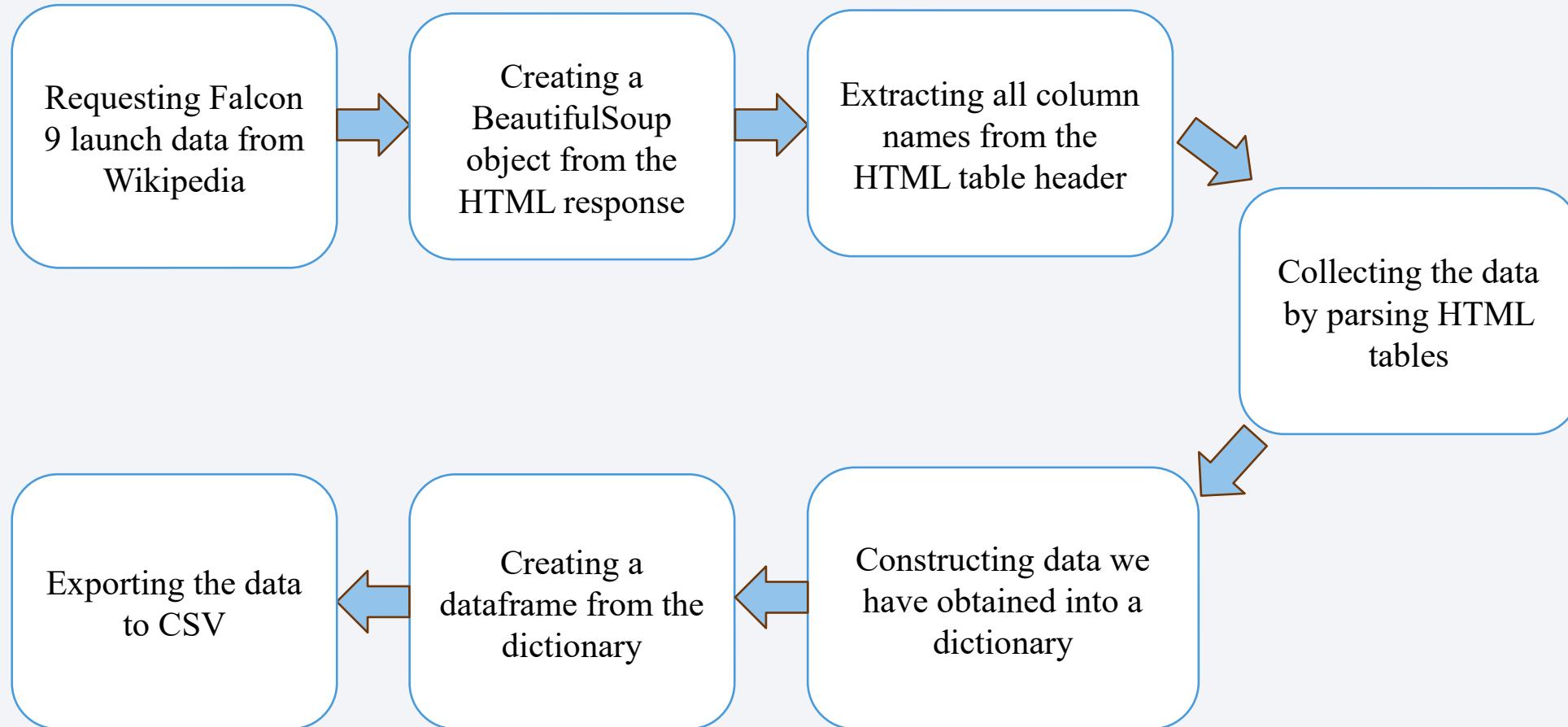
Exporting the data to CSV

Replacing missing values of Payload Mass column with calculated `.mean()` for this column

Filtering the dataframe to only include Falcon 9 launches

Creating a dataframe from the dictionary

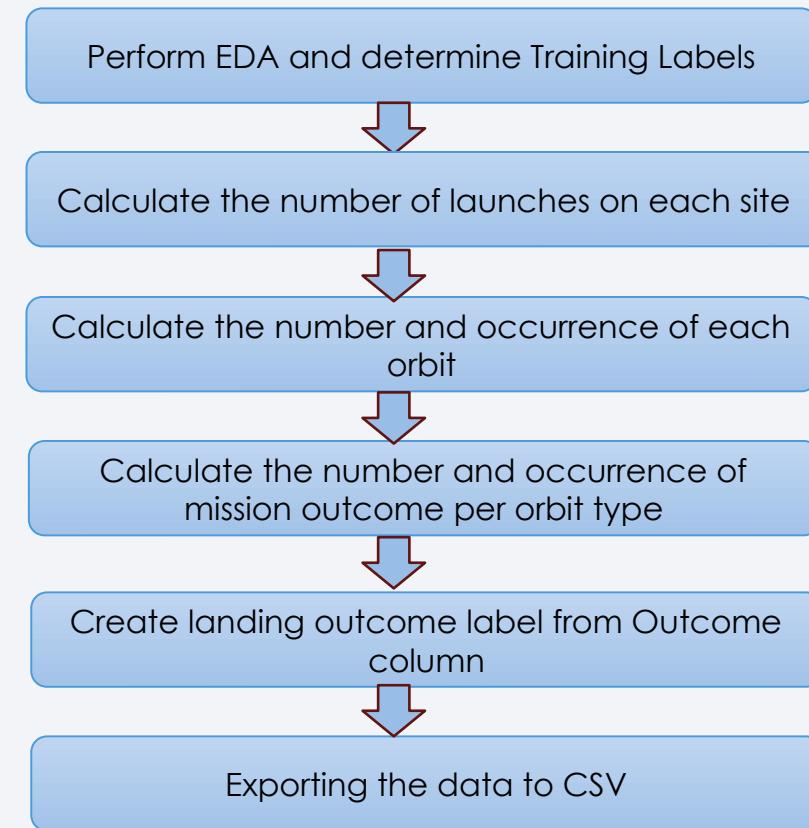
# Data Collection - Scraping



[GitHub URL: Webscraping](#)

# 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.



# EDA with Data Visualization

## Charts 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 and 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)

# 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

# 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 URL: 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:

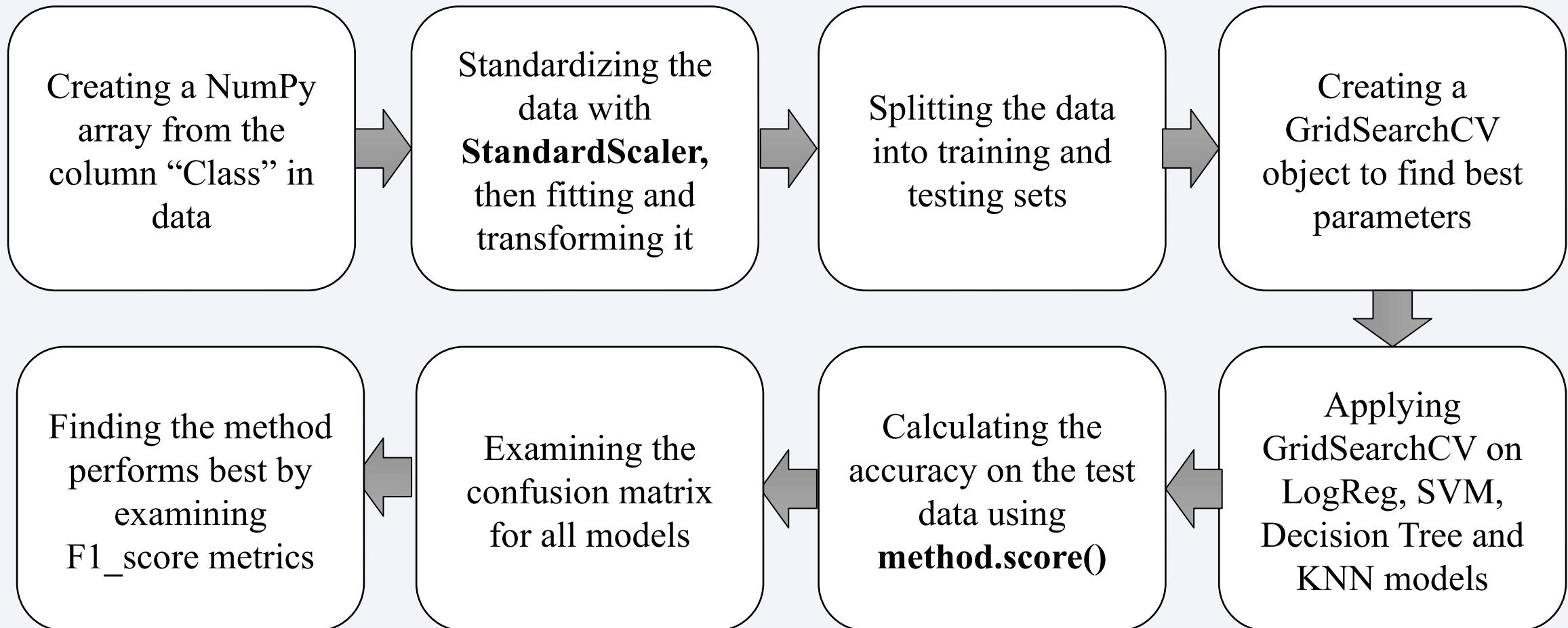
- Added a slider to select Payload range.

## 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 URL: Dashboard with Plotly](#)

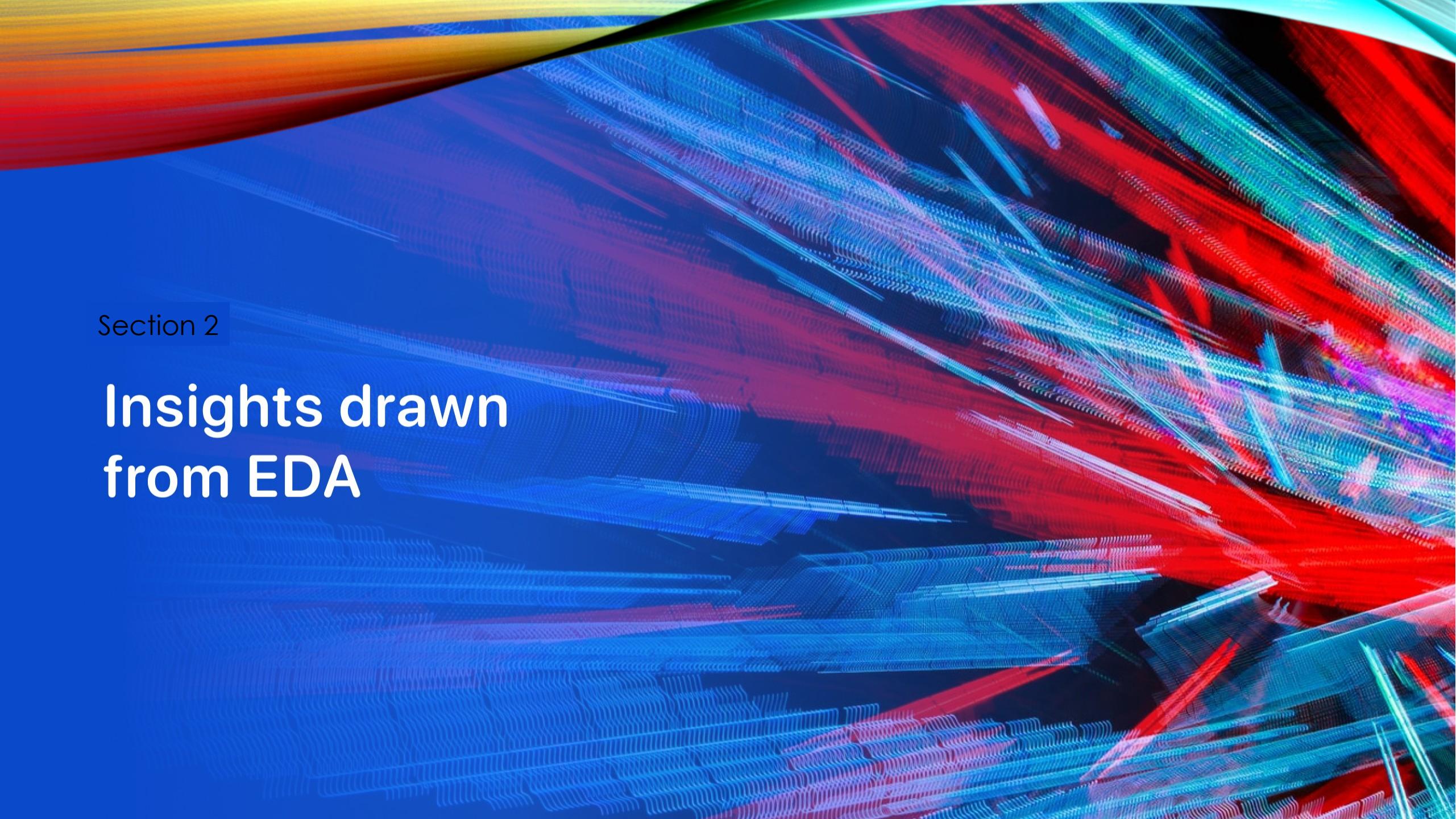
# Predictive Analysis (Classification)



## Results

---

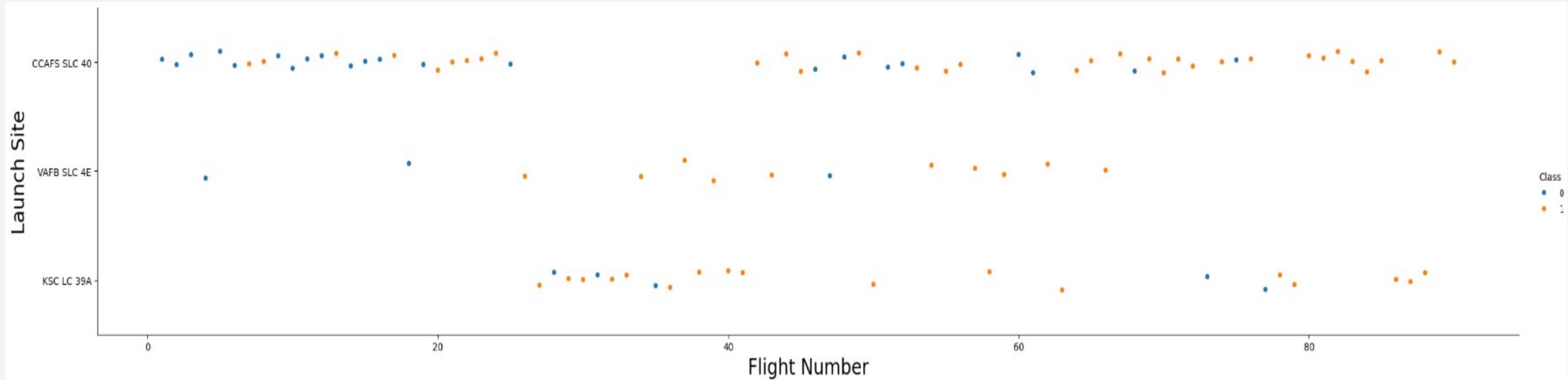
- Exploratory data analysis results
- Interactive analytics demo in screenshots
- Predictive analysis results



Section 2

## Insights drawn from EDA

## Flight Number vs. Launch Site



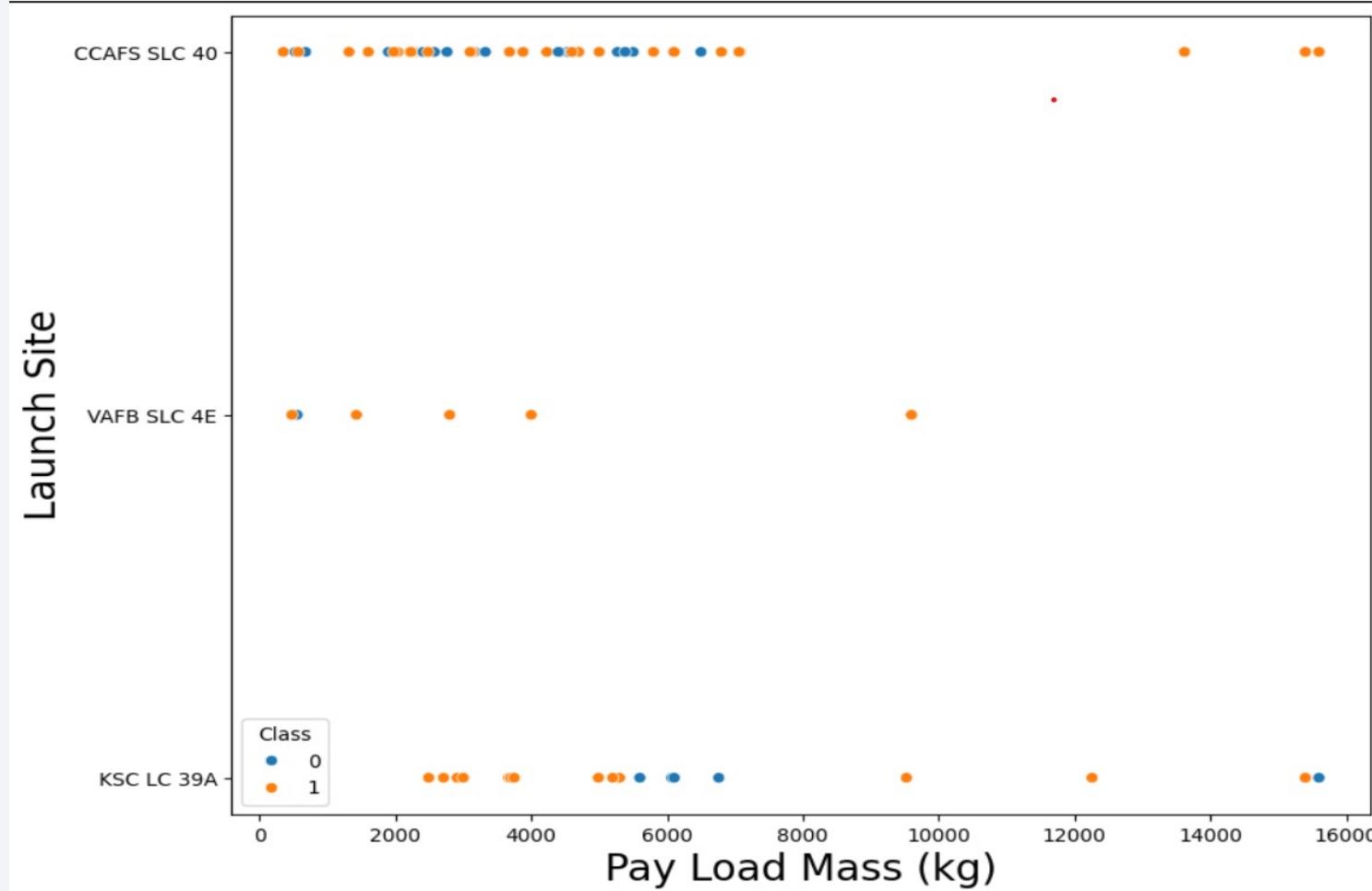
### Explanation:

- 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:

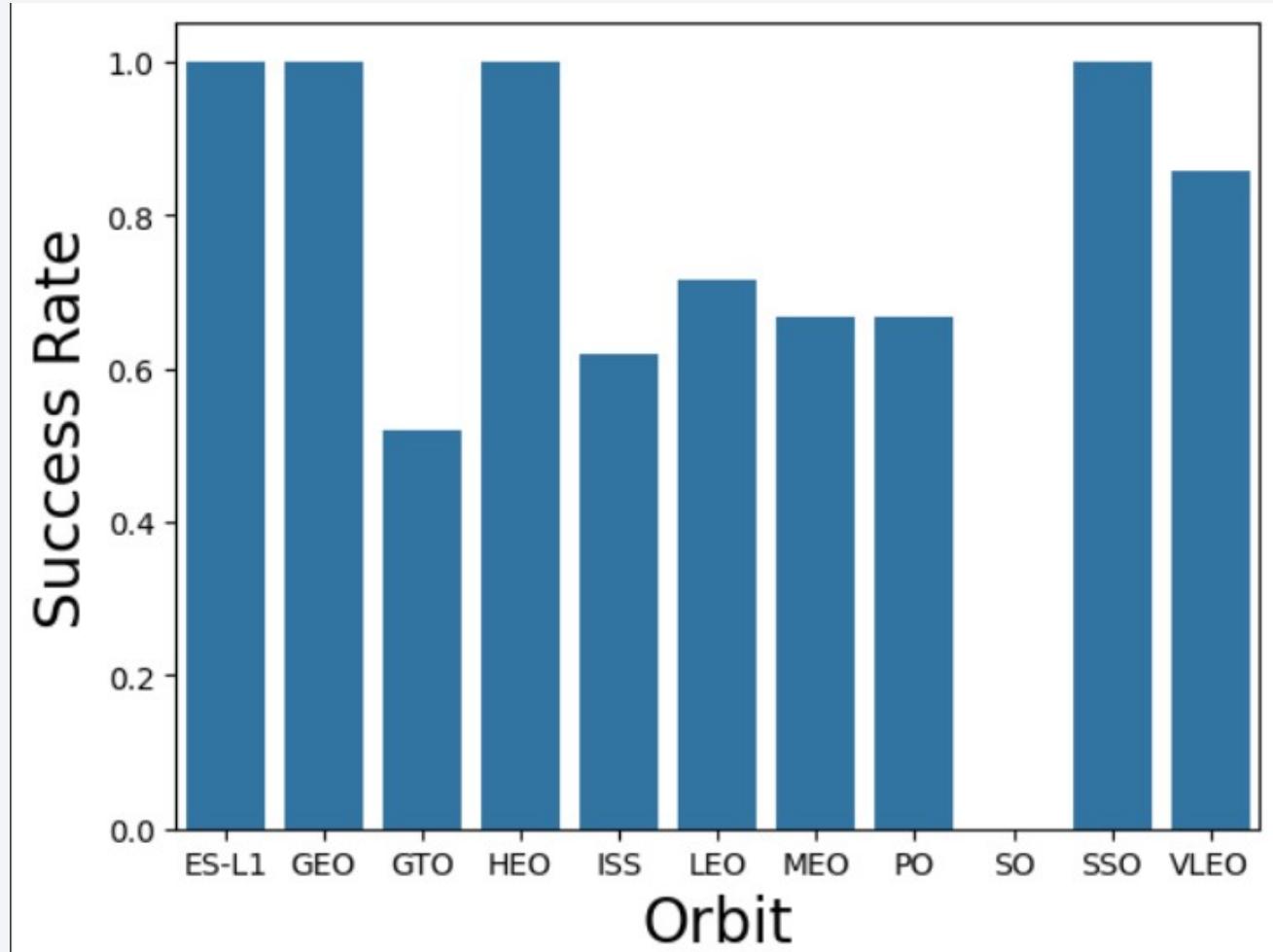
- For every launch site the higher the payload mass, the higher the success rate.
- Most of the launches with payload mass over 7000 kg were successful.
- KSC LC 39A has a 100% success rate for payload mass under 5500 kg too.



# Success Rate vs. Orbit Type

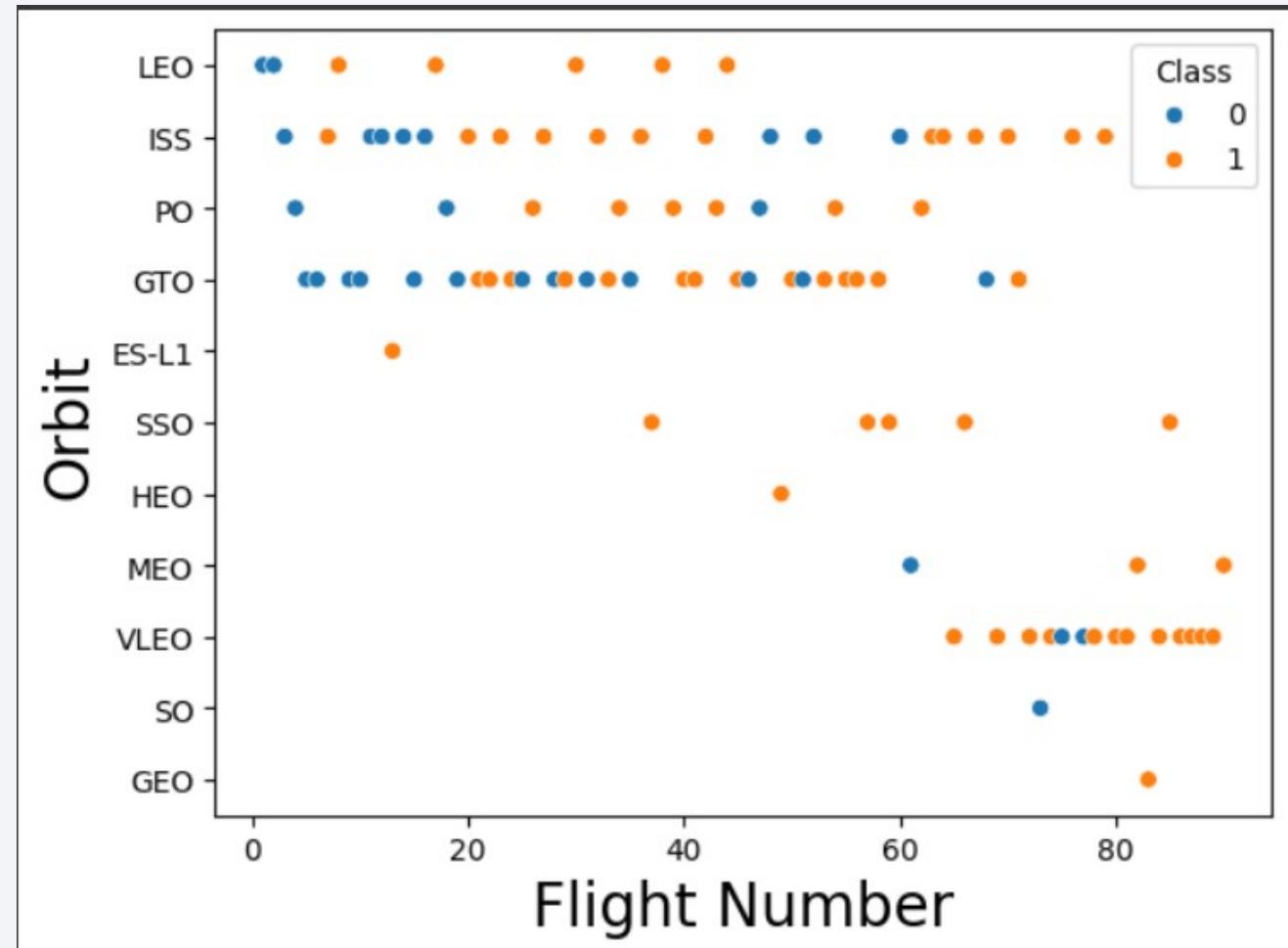
## Explanation:

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



# Flight Number vs. Orbit Type

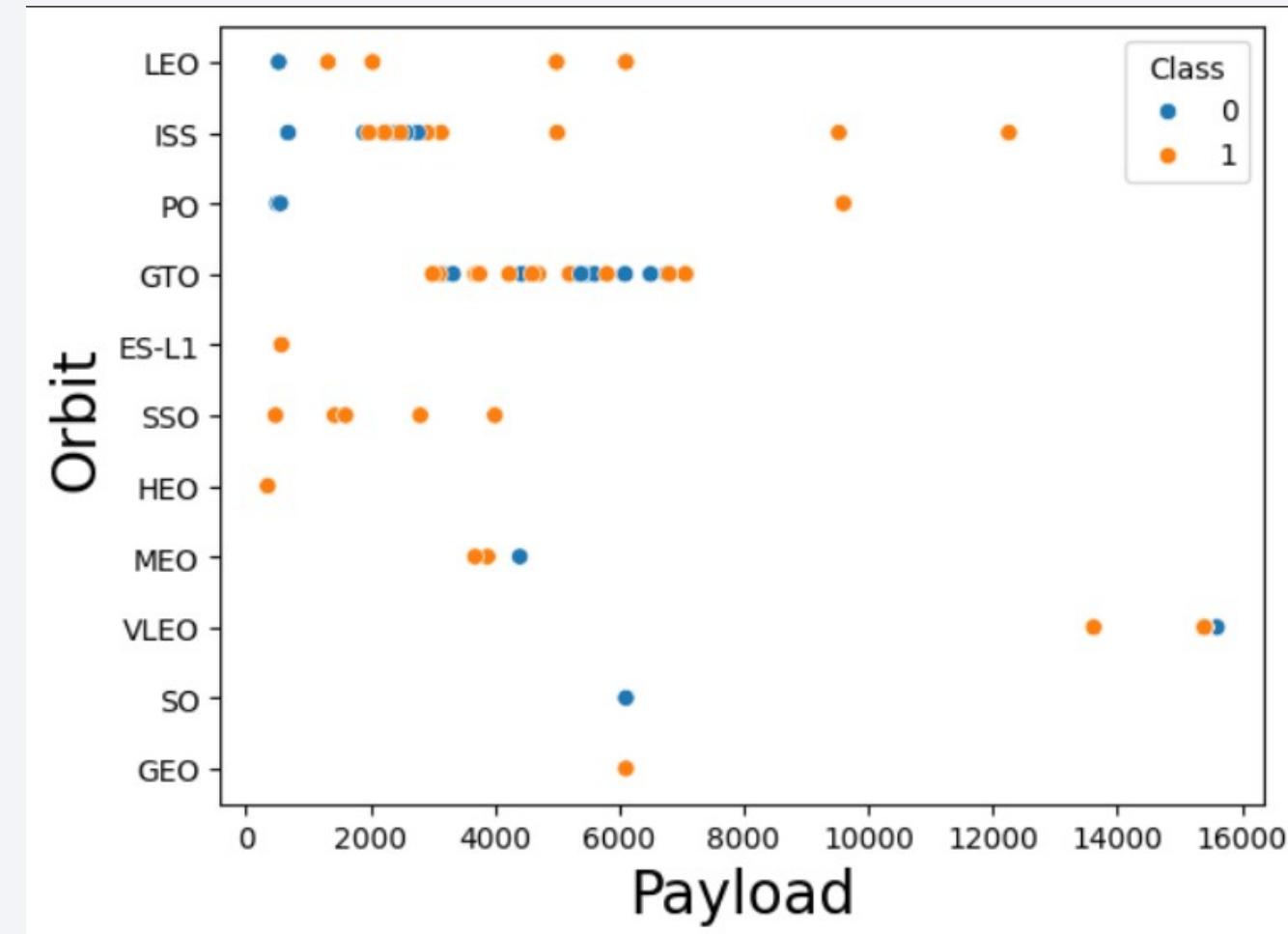
- **Explanation:**
- In the LEO orbit the Success appears related to the number of flights; on the other hand, there seems to be no relationship between flight number when in GTO orbit.



# Payload vs. Orbit Type

## Explanation:

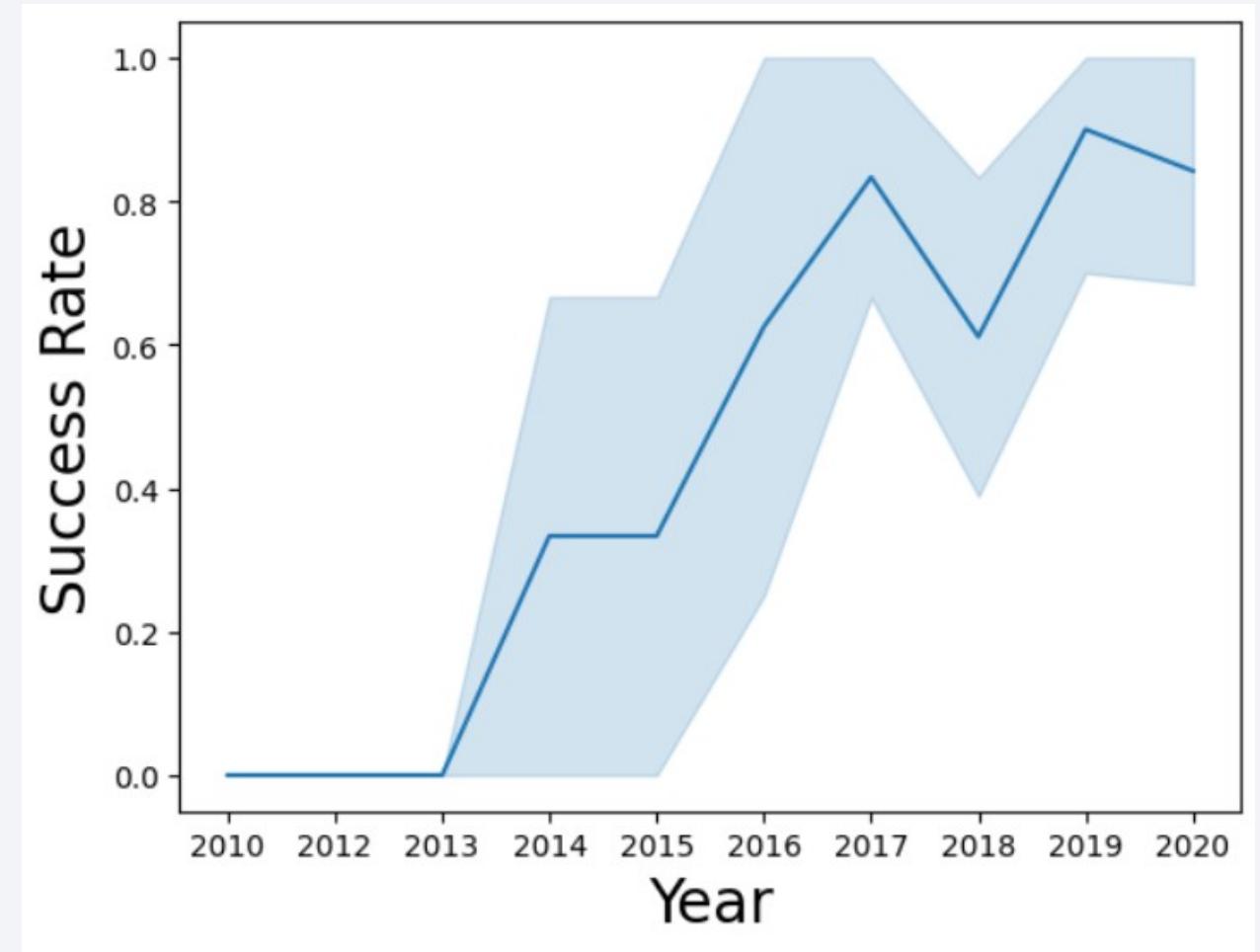
- Heavy payloads have a negative influence on GTO orbits and positive on GTO and Polar LEO (ISS) orbits



## Launch Success Yearly Trend

### Explanation:

- The success rate since 2013 kept increasing till 2020.



## EDA with SQL

## All Launch Site Names

```
▶ %sql select distinct Launch_Site from SPACEXTBL;
```

```
→ * sqlite:///my_data1.db
```

```
Done.
```

### 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 [5]: %sql select \* from SPACEXDATASET where launch\_site like 'CCA%' limit 5;

```
* ibm_db_sa://wzf08322:***@0c77d6f2-5da9-48a9-81f8-86b520b87518.bs2io90108kqb1od81cg.databases.appdomain.cloud:31198/bludb
Done.
```

Out[5]:

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	07:44:00	F9 v1.0 B0005	CCAFS LC-40	Dragon demo flight C2	525	LEO (ISS)	NASA (COTS)	Success	No attempt
2012-10-08	00:35:00	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

```
[12] %sql select sum(payload_mass_kg_) from spacextbl where customer = 'NASA (CRS)';  
→ * sqlite:///my_data1.db  
Done.  
sum(payload_mass_kg_)  
45596
```

## Average Payload Mass by F9 v1.1

```
▶ %sql select avg(payload_mass_kg_) from spacextbl where booster_version like 'F9 v1.1%'  
→ * sqlite:///my_data1.db  
Done.  
avg(payload_mass_kg_)  
2534.6666666666665
```

## First Successful Ground Landing Date

```
[14] %sql select min(Date) from spacextbl where Landing_Outcome = 'success (ground pad)';

→ * sqlite:///my_data1.db
Done.

min(Date)
2015-12-22
```

## Successful Drone Ship Landing with Payload between 4000 and 6000

```
[15] %sql select booster_version from spacextbl where Landing_Outcome='Success (drone ship)' and payload_mass_kg_ between 4000 and 6000;  
→ * 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

```
%sql select mission_outcome, count(*) as total_number from spacextbl group by mission_outcome;  
* sqlite:///my_data1.db  
Done.  


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


```

## Boosters Carried Maximum Payload

```
[17] %sql select booster_version from spacextbl where payload_mass_kg_ = (select max(payload_mass_kg_) from spacextbl);  
→ * sqlite:///my_data1.db  
Done.  
Booster_Version  
F9 B5 B1048.4  
F9 B5 B1049.4  
F9 B5 B1051.3  
F9 B5 B1056.4  
F9 B5 B1048.5  
F9 B5 B1051.4  
F9 B5 B1049.5  
F9 B5 B1060.2  
F9 B5 B1058.3  
F9 B5 B1051.6  
F9 B5 B1060.3  
F9 B5 B1049.7
```

# 2015 Launch Records

```
[18] %sql select booster_version, launch_site, Landing_Outcome from spacextbl where Landing_Outcome='Failure (drone ship)' and substr(Date,0,5)='2015';
```

```
→ * sqlite:///my_data1.db
```

Done.

<b>Booster_Version</b>	<b>Launch_Site</b>	<b>Landing_Outcome</b>
F9 v1.1 B1012	CCAFS LC-40	Failure (drone ship)
F9 v1.1 B1015	CCAFS LC-40	Failure (drone ship)

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

```
[19] %sql select Landing_Outcome, count(*) as Count from spacextbl where Date between '2010-06-04' and '2017-03-20' group by Landing_Outcome order by Count desc;
```

```
→ * sqlite:///my_data1.db
Done.

  Landing_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
Precluded (drone ship) 1
```



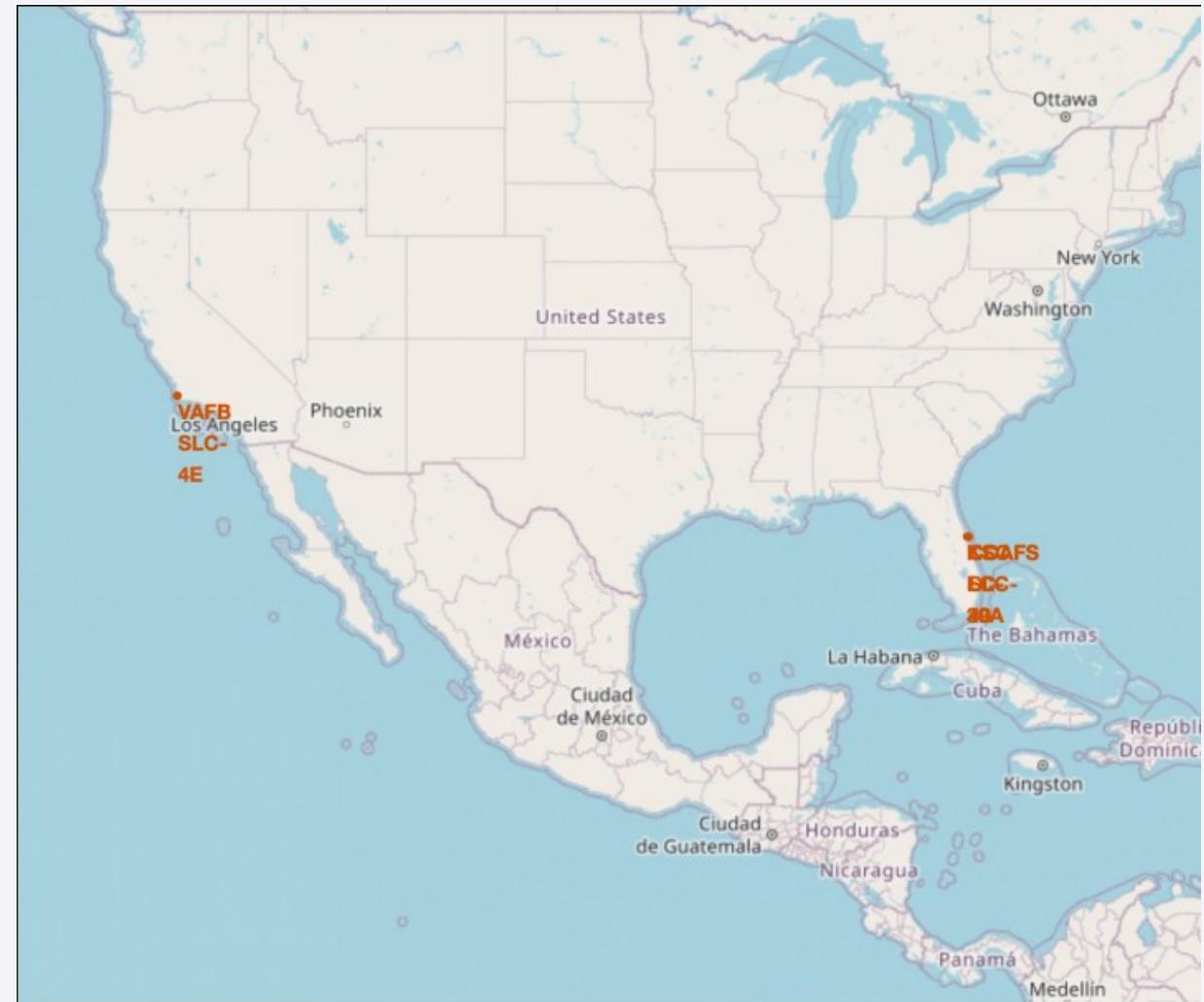
Section 3

# Launch Sites Proximities Analysis

# All launch sites' location markers on a global map

## Explanation:

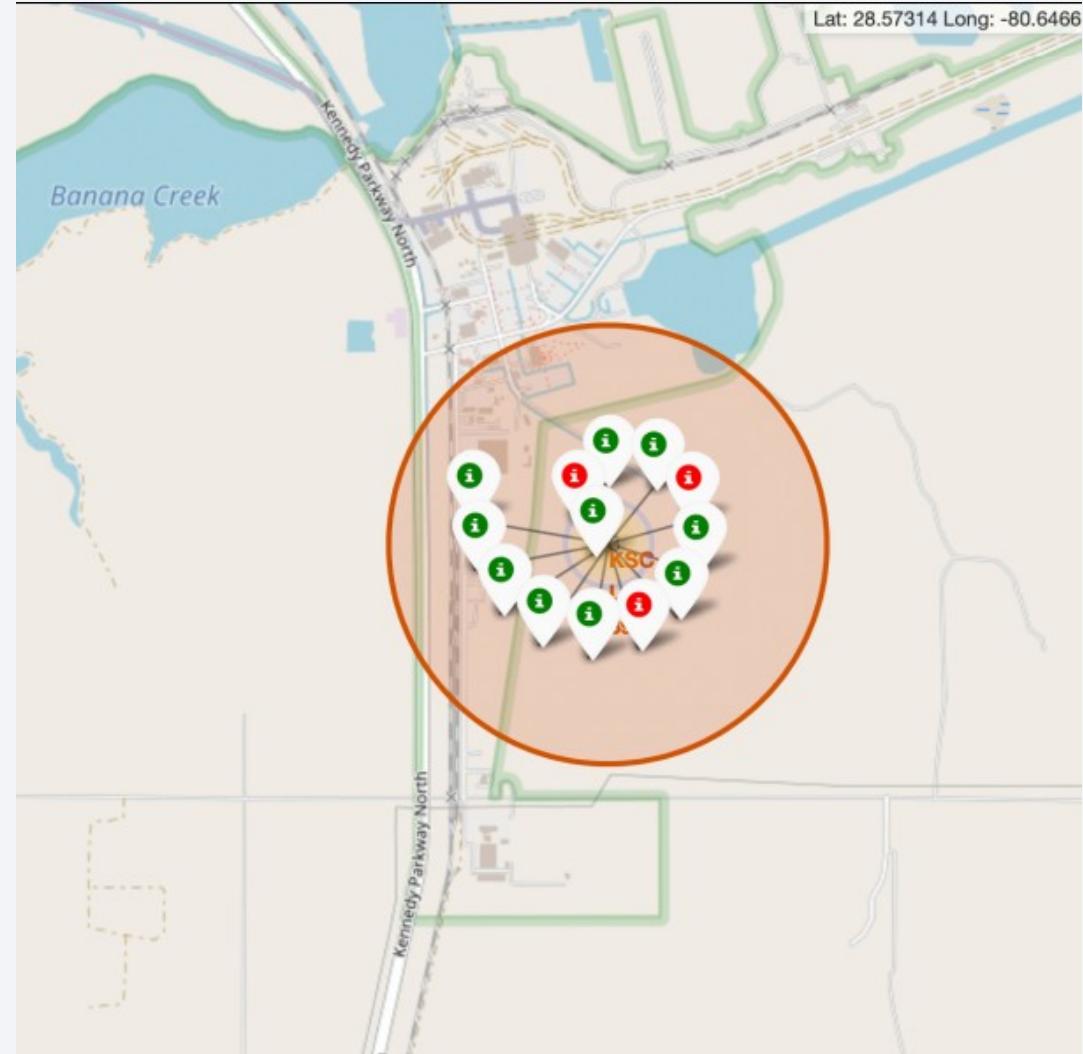
- Most of Launch sites are in proximity to the Equator line. The land is moving faster at the equator than any other place on the surface of the Earth. Anything on the surface of the Earth at the equator is already moving at 1670 km/hour. If a ship is launched from the equator it goes up into space, 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.



## Color-labeled launch records on the map

### Explanation:

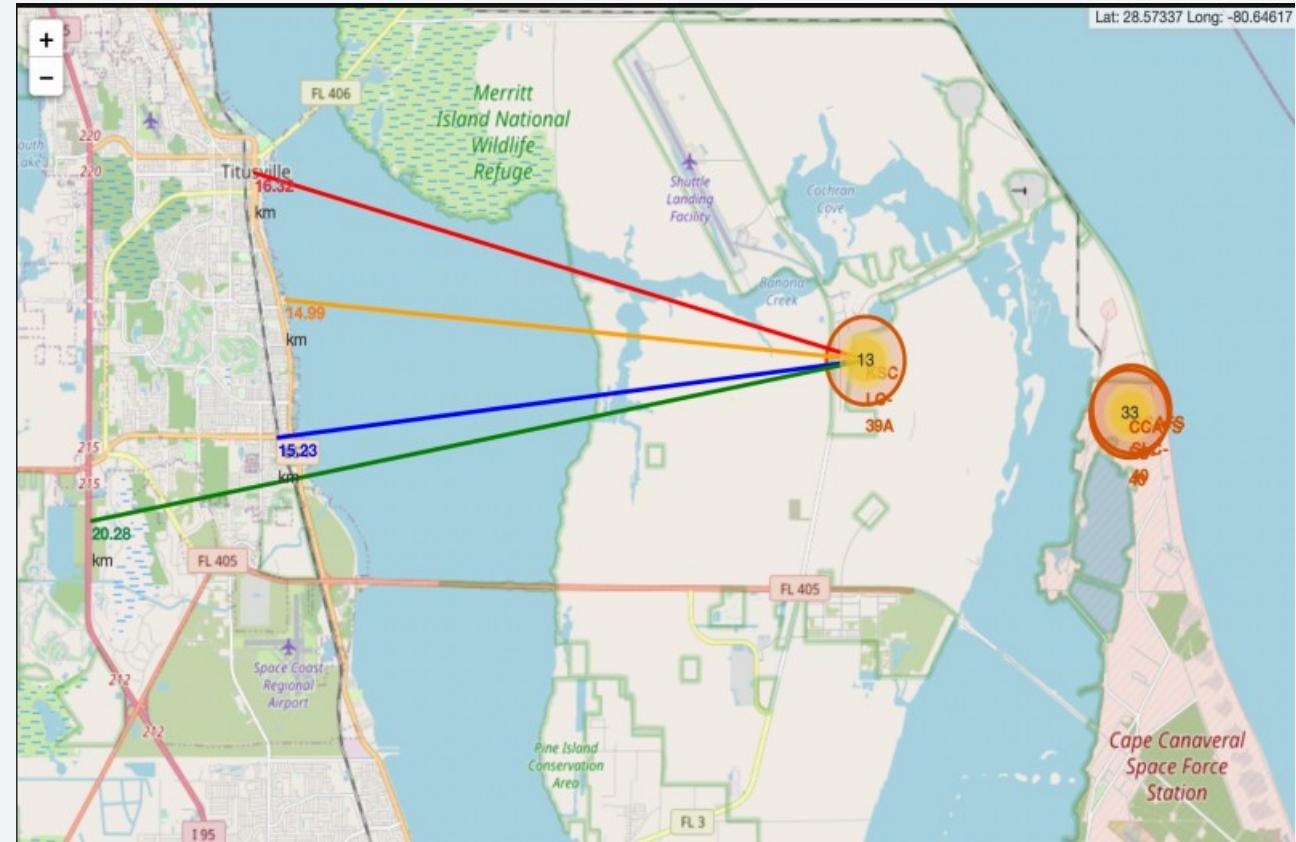
- From the colour-labeled markers we should be able to easily identify which launch sites have relatively high success rates.
- Green Marker = Successful Launch
- Red Marker = Failed Launch
- Launch Site KSC LC-39A has a very high Success Rate



# Distance from the launch site KSC LC-39A to its proximit

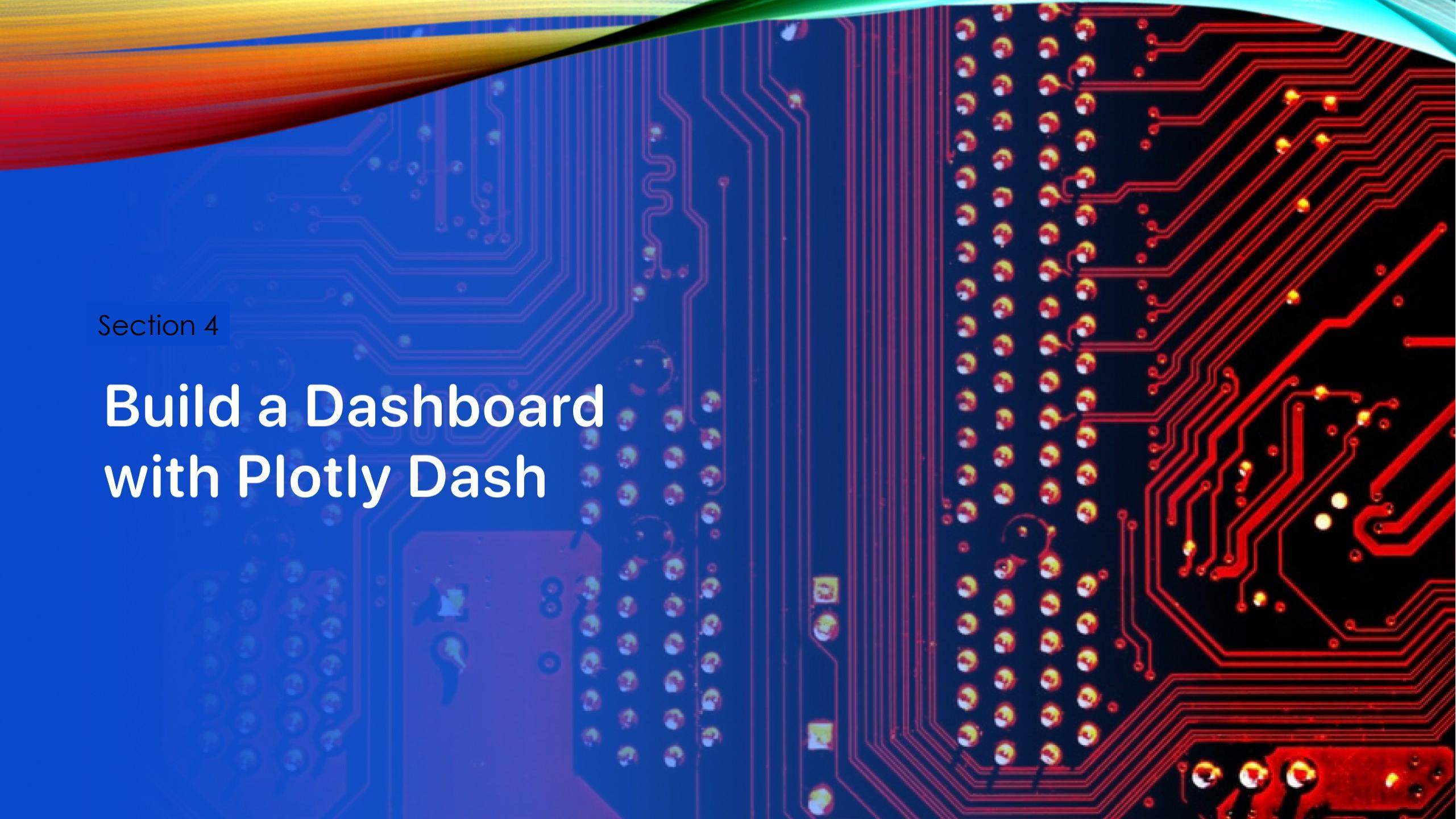
## Explanation:

- 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.

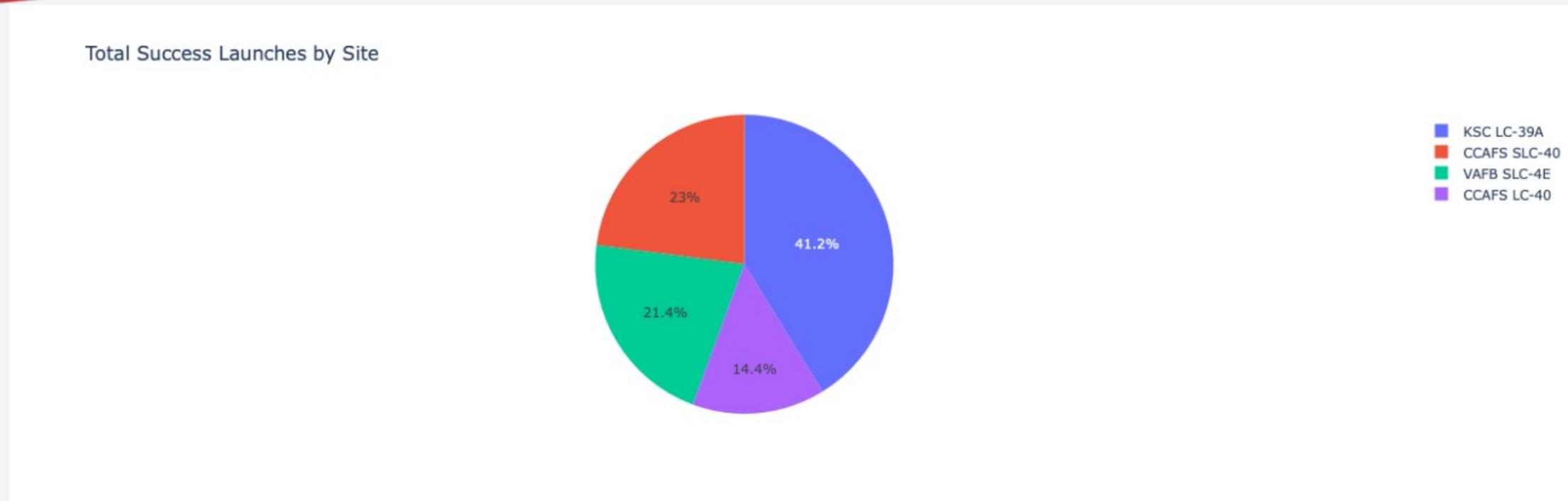


Section 4

# Build a Dashboard with Plotly Dash



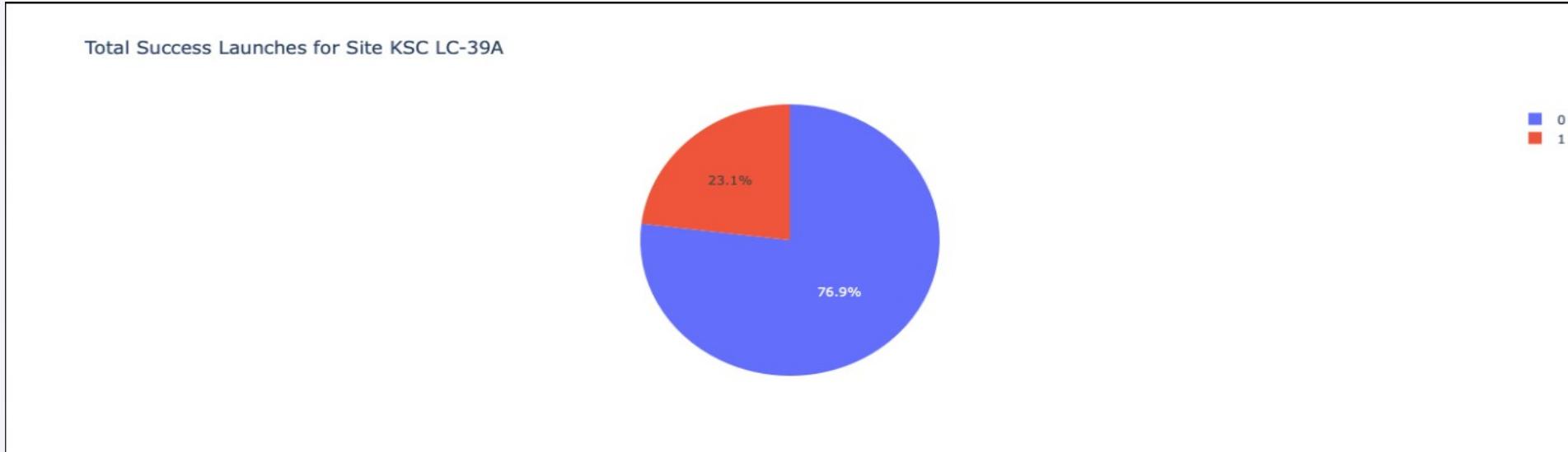
# Launch success count for all sites



## Explanation:

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

## Launch site with highest launch success ratio



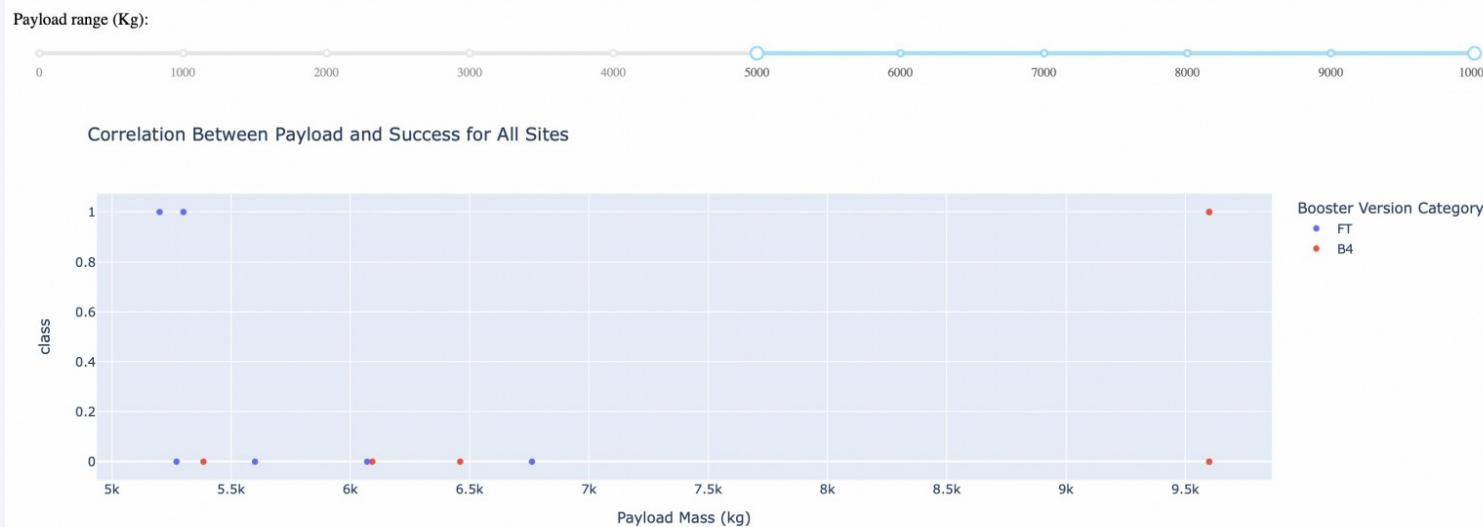
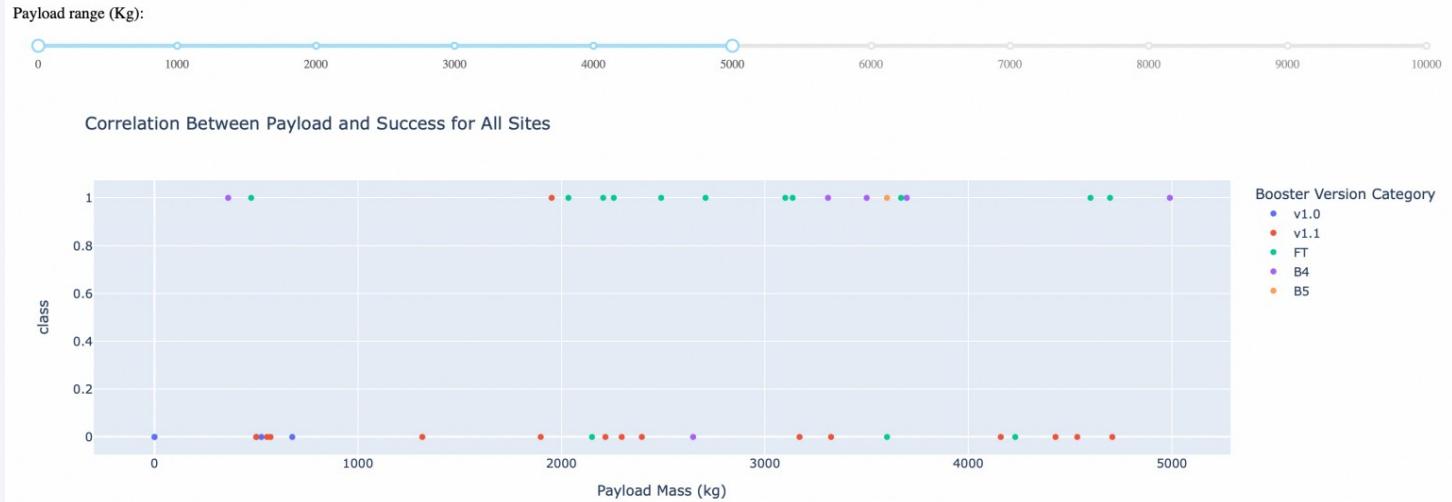
### Explanation:

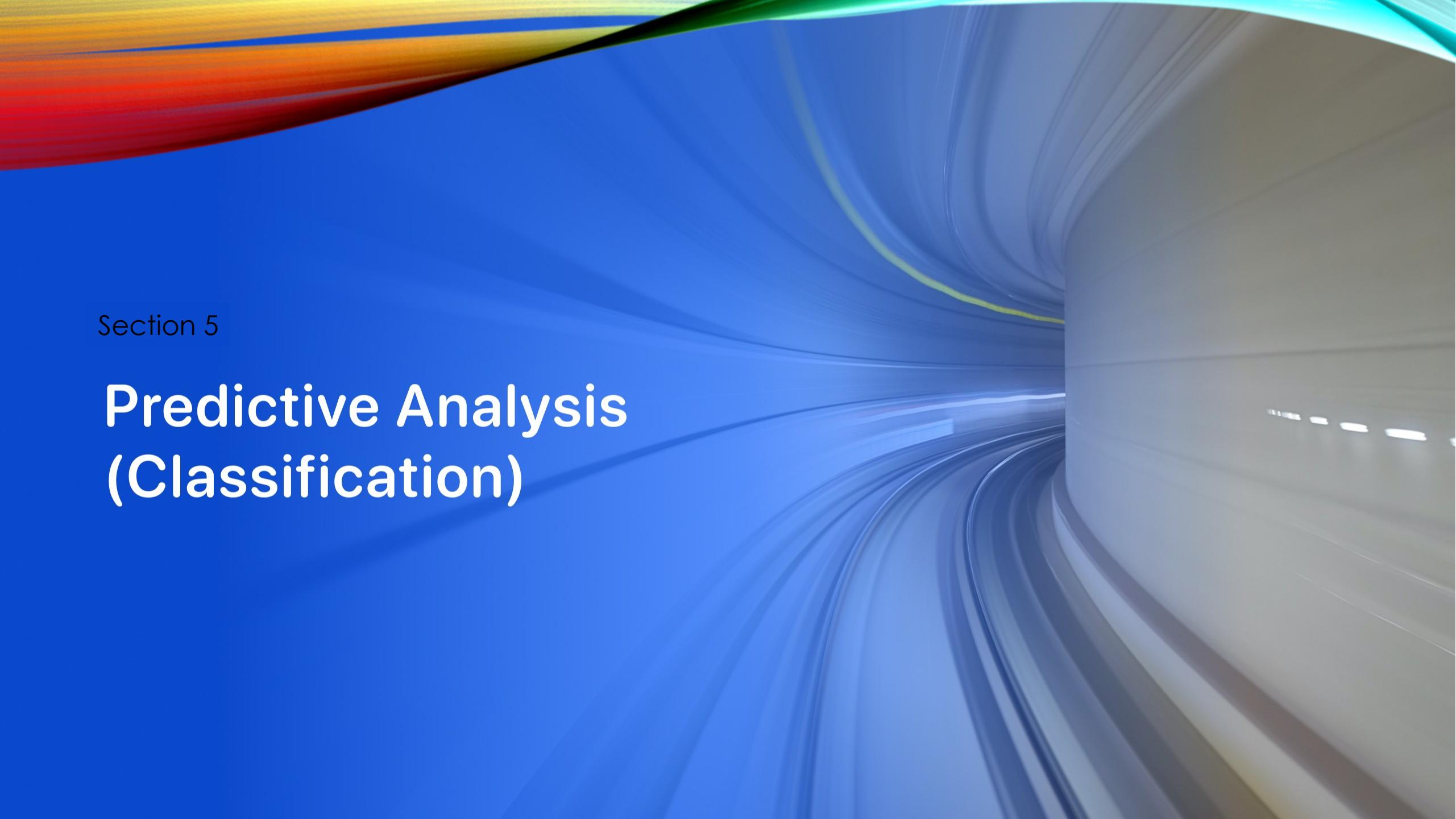
- KSC LC-39A has the highest launch success rate (76.9%) with 10 successful and only 3 failed landings.

# Payload Mass vs. Launch Outcome for all sites

## Explanation:

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



The background of the slide features a dynamic, abstract design. It consists of several curved, overlapping bands of color that create a sense of motion. The colors transition from warm tones like red, orange, and yellow at the top left to cooler tones like blue, green, and cyan towards the top right. Below these, there are darker, more saturated blue and purple bands that suggest a tunnel or a deep space. The overall effect is one of speed, technology, and data flow.

Section 5

# Predictive Analysis (Classification)

# Classification Accuracy

```
: print("Logistic Regression Test Accuracy :", logreg_test_score)
print("SVM Test Accuracy           :", svm_test_score)
print("Decision Tree Test Accuracy:", tree_test_score)
print("KNN Test Accuracy          :", knn_test_score)

accuracies = {
    'LogReg': logreg_test_score,
    'SVM': svm_test_score,
    'DecisionTree': tree_test_score,
    'KNN': knn_test_score
}
best_model = max(accuracies, key=accuracies.get)
print(f"\nBest performing model on test data is: {best_model} with accuracy {accuracies[best_model]:.2f}")
```

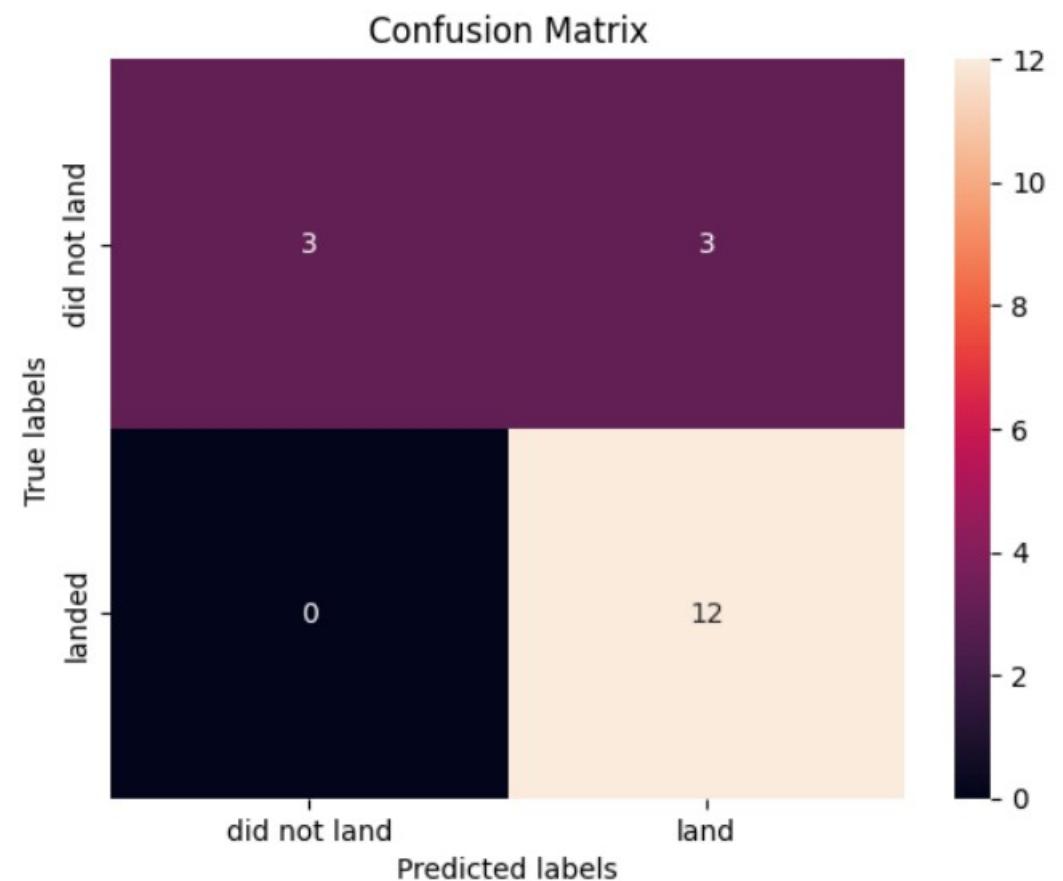
```
Logistic Regression Test Accuracy : 0.8333333333333334
SVM Test Accuracy           : 0.8333333333333334
Decision Tree Test Accuracy: 0.8333333333333334
KNN Test Accuracy          : 0.8333333333333334
```

```
Best performing model on test data is: LogReg with accuracy 0.83
```

# 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

- a. Decision Tree Model is the best algorithm for this dataset.
- b. 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.
- c. The success rate of launches increases over the years.
- d. KSC LC-39A has the highest success rate of the launches from all the sites.
- e. Orbits ES-L1, GEO, HEO and SSO have 100% success rate.





Thank you!