



IBM Developer  
SKILLS NETWORK

# Winning Space Race with Data Science

Felix Tandeas  
20 October 2023



# 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 methodologies

- Exploratory Data Analysis Results
- Interactive Analytics demo
- Predictive Analysis results

# Introduction

---

## Project background and context

SpaceX is one of the most successful companies of the commercial space age, making space travel affordable. SpaceX advertises Falcon 9 rocket launches on its website with a cost of 62 million dollars whereas 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 predict whether the first stage will launch, we can determine the cost of the launch. Based on public information and machine learning models we are going to predict if SpaceX will reuse its first stage.

## Problems you want to find answers

- How do variables such as payload mass, launch site, number of flights and orbits affect the success of the first stage landing?
- Does the rate of successful landings increase over the years?
- Which machine learning model can classify the landing outcome with highest accuracy?





Section 1

# Methodology

# Methodology

---

## Executive Summary

- Data collection methodology:
  - SpaceX Rest API
  - Web scrapping from Wikipedia
- Perform data wrangling
  - Filtering the data
  - Dealing with missing values
  - Using one hot encoding to prepare the data for 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 evaluating classification model for maximum accuracy

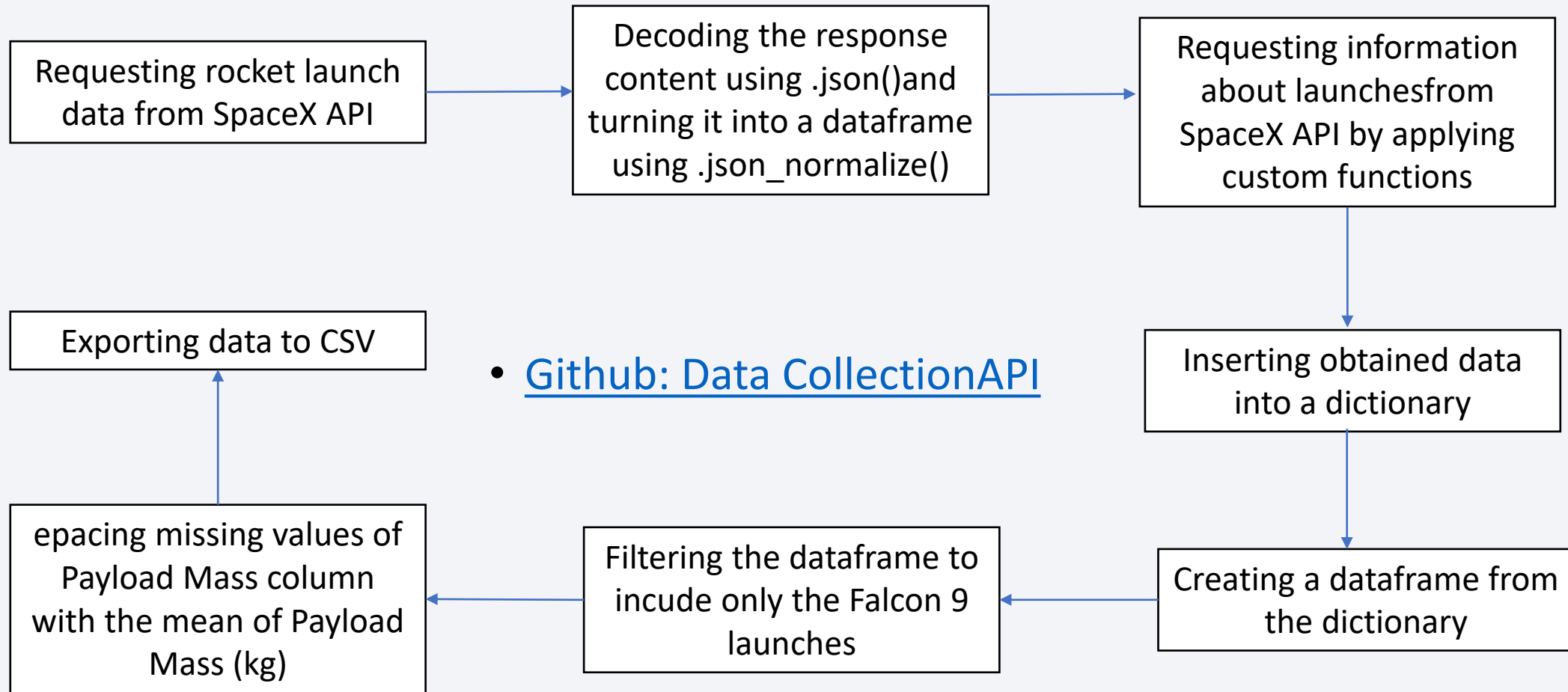
# Data Collection

---

- Data collection process involved a combination of API requests from SpaceX REST API and Web Scraping data from a table in SpaceX's wikipedia page
- Both methods had to be used to get complete information about the launches for an in depth analysis.
- Fields collected by using the SpaceX REST API:
  - FlightNumber ▪ Date ▪ BoosterVersion ▪ PayloadMass ▪ Orbit ▪ LaunchSite ▪ Outcome ▪ Flights ▪ GridFins ▪ Reused ▪ Legs ▪ LandingPad ▪ Block ▪ ReusedCount ▪ Serial ▪ Longitude ▪ Latitude
- Fields collected by using Wikipedia web scrapping:
  - FlightNo ▪ Launch Site ▪ Payload ▪ PayloadMass ▪ Orbit ▪ Customer ▪ Launch outcome ▪ Version  
Booster ▪ Booster Landing ▪ Date ▪ Time

# Data Collection – SpaceX API

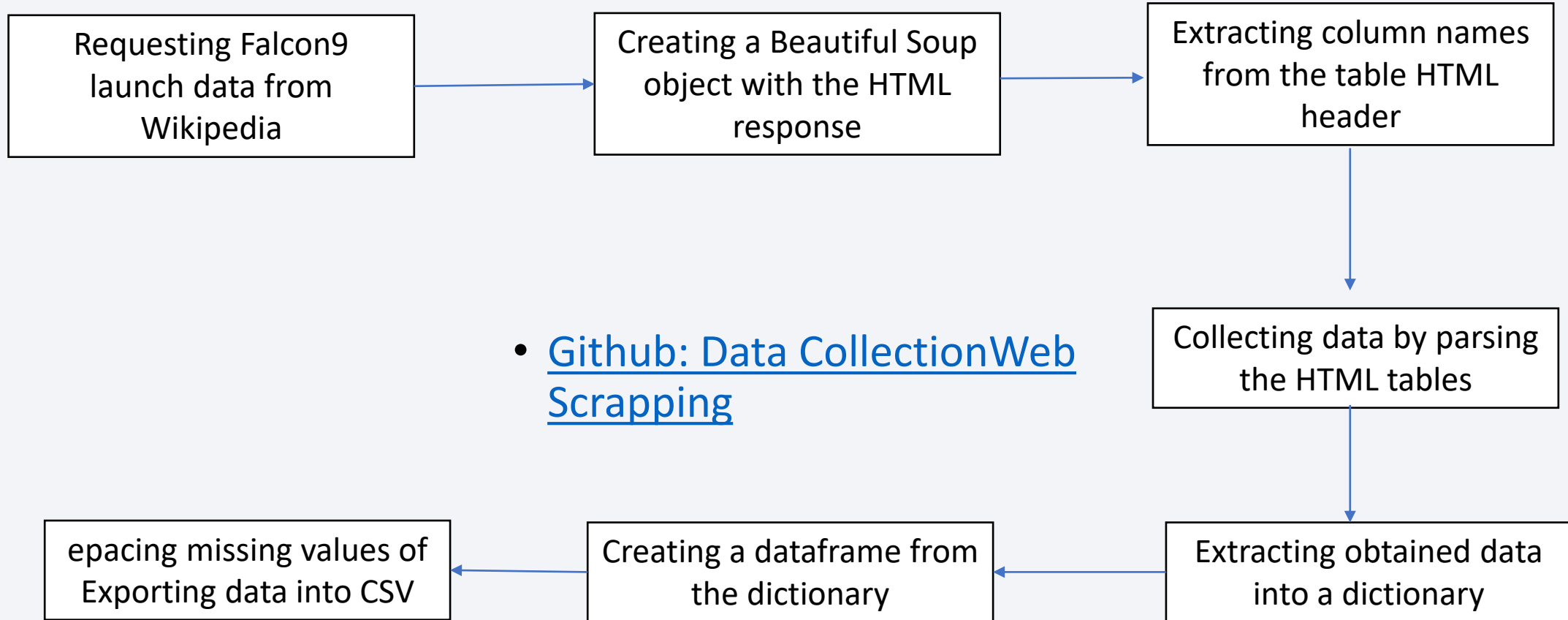
---





# Data Collection - 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, 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.

Those outcomes are mainly converted into Training Labels with 1 means the booster successfully landed, 0 means it was unsuccessful.

- [Github: Data Wrangling](#)

Calculate number of launches for each site

Calculate the number and occurrence of each orbit

Calculate the number and occurrence and mission outcome per orbit type

Create a landing outcome label from Outcome column

Exporting data to CSV

# EDA with Data Visualization

---

## Chart Plotted:

- Flight Number vs Payload Mass
  - Flight Number vs Launch Site
  - Payload Mass vs Launch Site
  - Success rate per Orbit Type
  - Flight Number vs Orbit Type
  - Payload Mass vs Orbit Type
  - Yearly Success Rate
- Scatter plots are used to show to relationship between different variables. If there is a relationship, they could be used in machine learning model for classifying landing outcomes.
  - Bar chart is used to show the contribution of categorical variables to the landing success rate.
  - Line chart is used to show how the success rate changes over time.
- [Github: Data Visualitation](#)

# EDA with SQL

---

## 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 per month for the 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: SQL](#)

# Build an Interactive Map with Folium

---

## Markers of all launch sites

- Added marker with Circle, Popup Label and Text Label for NASA Johnson Space Center using its latitude and longitude coordinates as a start location
- Added marker with Circle, Popup Label and Text Label for all launch sites using their latitude and longitude coordinates to show their geographical location and proximity to Equator and coasts

## Colored markers of launch outcomes for each launch site

- Added colored markers of successful (green) and failed (red) launch outcomes per launch site to display launch sites with highest success rates

## Distances between a launch site to its proximities

- Added colored lines to show the distance between launch site KSC LC-39A and its proximities including Railway, Highway, Coastline and closest city.

- [Github: Folium Map](#)



# Build a Dashboard with Plotly Dash

---

- Launch Sites Dropdown List :

Added a dropdown list to enable launch site selection

- Pie Chart showing success launches for all sites:

Added a pie chart to display the total count of successful launches for all sites and the Success vs Failed launches count for a site, if a specific site was selected

- Slider of Payload Mass range:

Added a pie chart to display the total count of successful launches for all sites and the Success vs Failed launches count for a site, if a specific site was selected.

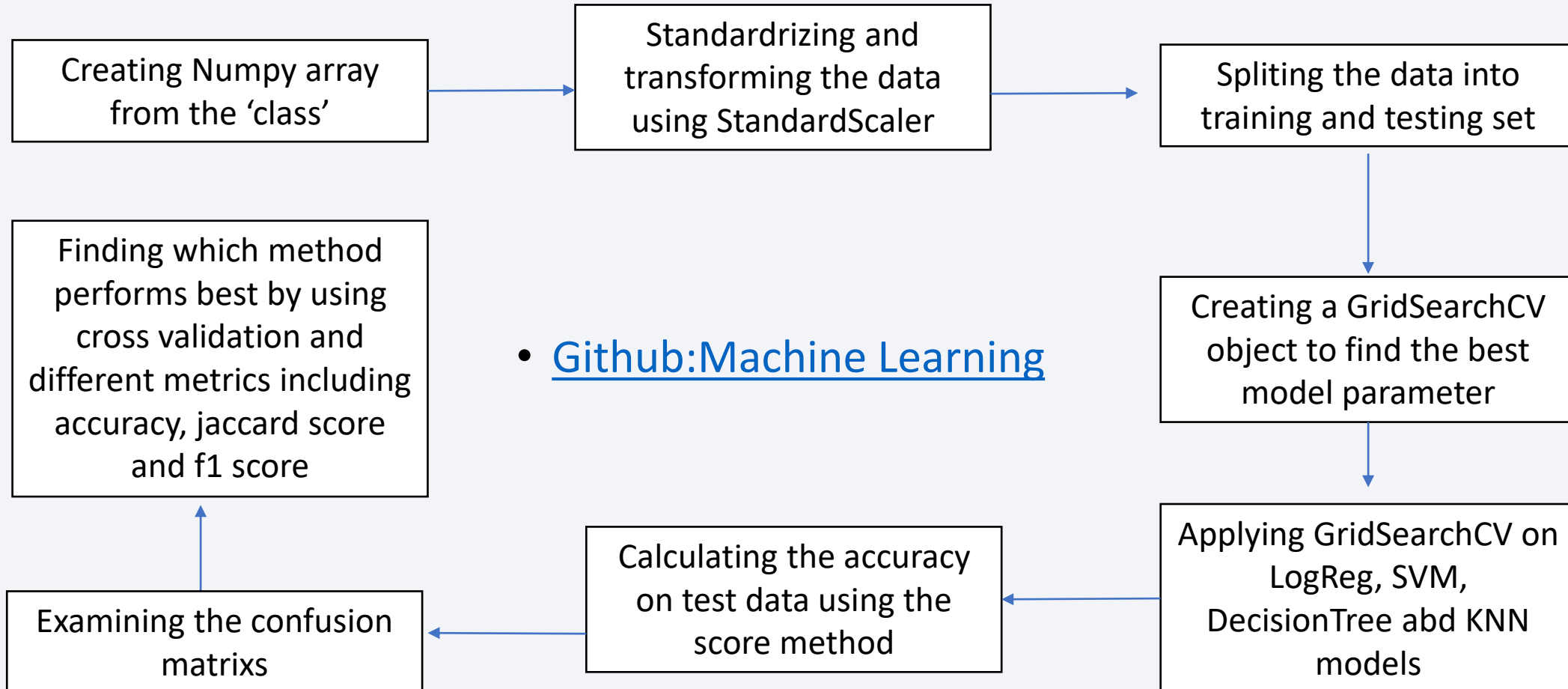
- Scatter plot of Payload Mass vs Success Rate for different Booster Versions:

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

- [Github: Dash App](#)

# Predictive Analysis (Classification)

---



# Results

---

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



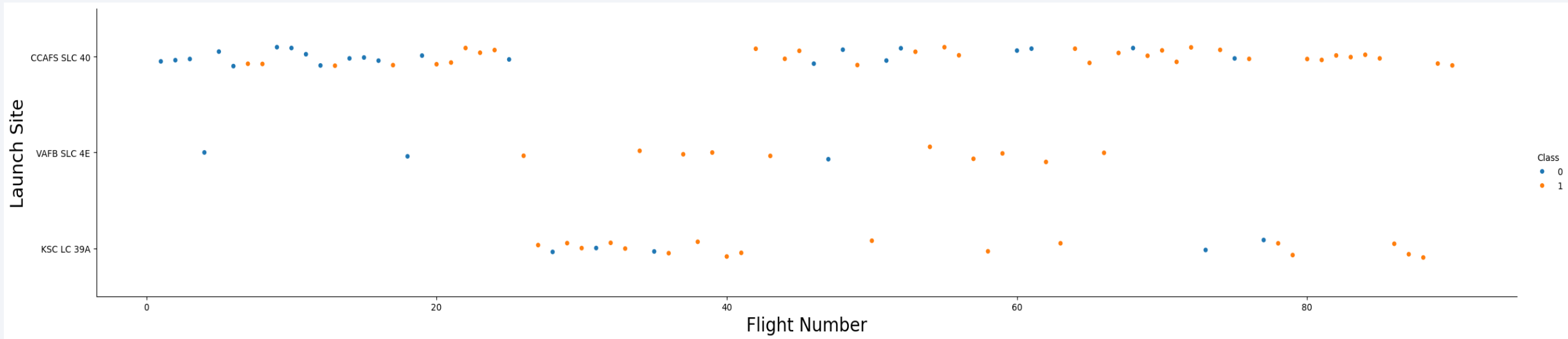
The background of the slide is an abstract composition. It features a dark blue base color. Overlaid on this are numerous diagonal streaks in shades of red and cyan. A faint, light blue grid pattern is also visible, particularly in the lower half of the image. The overall effect is dynamic and technological.

Section 2

# Insights drawn from EDA



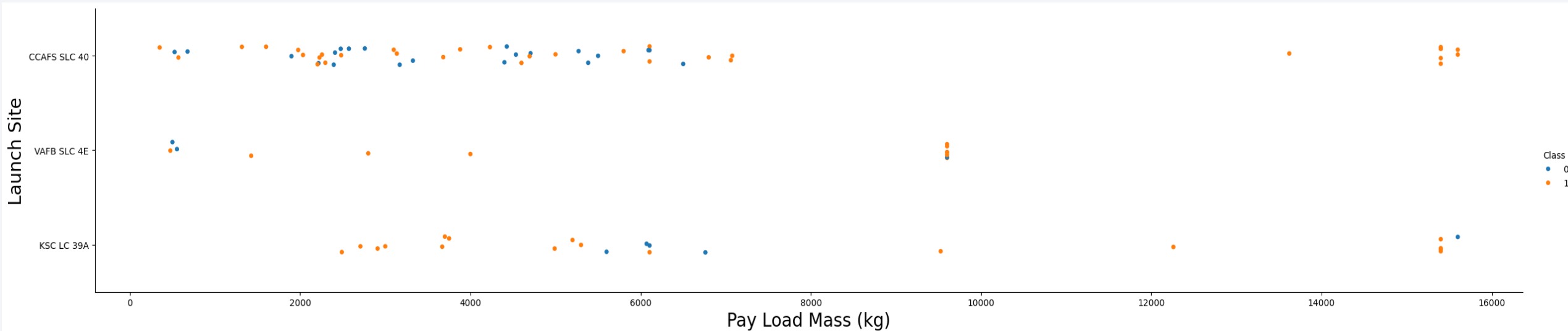
# Flight Number vs. Launch Site



- The earliest flights failed while the latest flights succeeded
- The launch site CCAFS SCL 40 has about half of all launches
- VAFB SLC 4E and KSC LC 39A have higher success rates

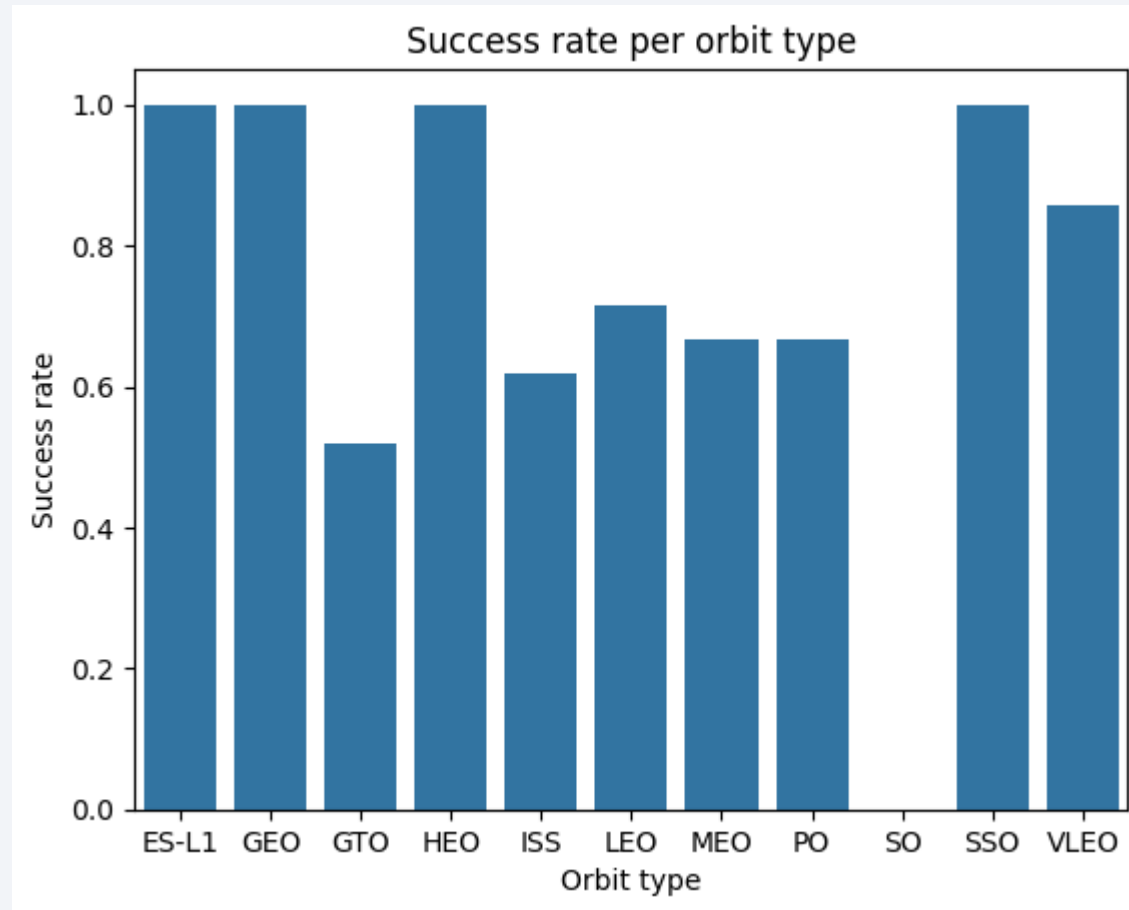


# Payload vs. Launch Site



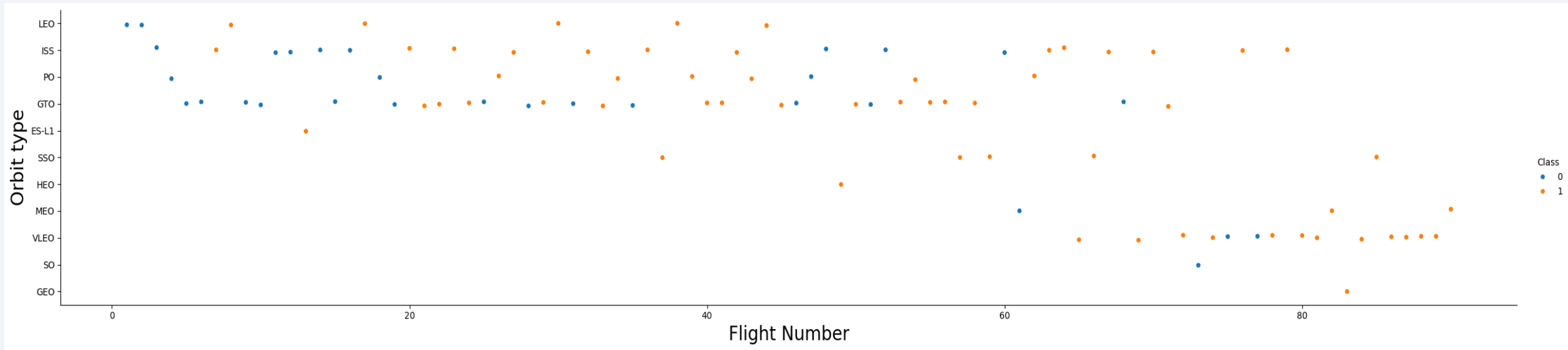
- The higher the payload mass, the higher the success rate for each launch site
- Most of the launches with payload mass over 7000 kg were successful
- KSC LC 39A has 100% success rate for payload mass under 5500 kg

# Success Rate vs. Orbit Type



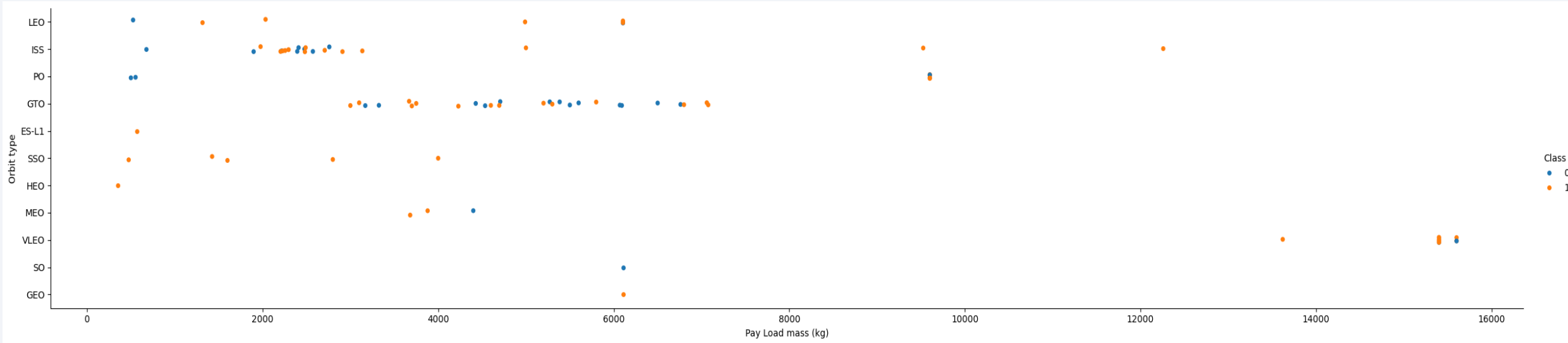
- Orbits ES-L1,GEO,HEO AND SSO have 100% success rate
- Orbits GTO,ISS,LEO,MEO and PO have success rate between 50% and 80%
- Orbit SO has 0% success rate

# Flight Number vs. Orbit Type



In the LEO orbit the success rate 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



- With heavy payloads the successful landing or positive landing rate are more for Polar, LEO and ISS.
- However for GTO we cannot distinguish this well as both positive landing rate and negative landing (unsuccessful mission) are here and there

# Launch Success Yearly Trend

---



- Since 2013 till 2020 success rate kept increasing



# All Launch Site Names

---

## Task 1

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

```
%%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

- Launch sites in the space mission

# Launch Site Names Begin with 'CCA'

## Task 2

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

[9]:

```
%%sql
```

```
select*from SPACEXTABLE
where "Launch_Site" like 'CCA%' limit 5
```

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

Done.

[9]:

Date	Time (UTC)	Booster_Version	Launch_Site	Payload	PAYLOAD_MASS_KG_	Orbit	Customer	Mission_Outcome	Landing_Outcome
2010-04-06	18:45:00	F9 v1.0 B0003	CCAFS LC-40	Dragon Spacecraft Qualification Unit	0	LEO	SpaceX	Success	Failure (parachute)
2010-08-12	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
2012-08-10	00:35:00	F9 v1.0 B0006	CCAFS LC-40	SpaceX CRS-1	500	LEO (ISS)	NASA (CRS)	Success	No attempt
2013-01-03	15:10:00	F9 v1.0 B0007	CCAFS LC-40	SpaceX CRS-2	677	LEO (ISS)	NASA (CRS)	Success	No attempt

- 5 records where launch sites begin with 'CCA'

# Total Payload Mass

## Task 3

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

In [10]:

```
%%sql

select Customer, sum("PAYLOAD_MASS__KG_") as Total_payload_mass
from SPACEXTABLE
where Customer = 'NASA (CRS)'
```

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

Out[10]:

Customer	Total_payload_mass
NASA (CRS)	45596

- Total Payload Mass carried by booster launched by NASA(CRS) : 45596

# Average Payload Mass by F9 v1.1

---

## Task 4

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

```
In [11]: %%sql
select "Booster_version", avg("PAYLOAD_MASS__KG_") as Average_payload_mass
from SPACEXTABLE
where "Booster_version" = 'F9 v1.1'
```

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

```
Out[11]: 




```

- Average Payload Mass carried by booster version F9 v1.1. : 2928.4

# First Successful Ground Landing Date

---

## Task 5

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

*Hint: Use min function*

In [12]:

```
%%sql  
  
select "Landing_Outcome", min(Date) as First_successful_landing  
from SPACEXTABLE  
where "Landing_Outcome" = 'Success (ground pad)'
```

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

Done.

Out[12]:

Landing_Outcome	First_successful_landing
Success (ground pad)	2015-12-22

- The first successful landing outcome on ground pad on 22 Dec 2015



# Successful Drone Ship Landing with Payload between 4000 and 6000

## 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

```
[3]: %%sql

select distinct("Booster_Version") as Booster_Version_successful_in_drone_ship
from SPACEXTABLE
where "Landing_Outcome" = 'Success (drone ship)' and "PAYLOAD_MASS_KG_" between 4000 and 6000
```

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

Done.

```
[3]: Booster_Version_successful_in_drone_ship
```

F9 FT B1022

F9 FT B1026

F9 FT B1021.2

F9 FT B1031.2

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

# Total Number of Successful and Failure Mission Outcomes

## Task 7

List the total number of successful and failure mission outcomes

In [14]:

```
%%sql  
  
select "Mission_Outcome", count("Mission_Outcome")  
from SPACEXTABLE  
group by "Mission_Outcome"
```

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

Done.

Out[14]:

Mission_Outcome	count("Mission_Outcome")
Failure (in flight)	1
Success	98
Success	1
Success (payload status unclear)	1

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

# Boosters Carried Maximum Payload

## Task 8

List the names of the booster\_versions which have carried the maximum payload mass. Use a subquery

```
15]: %%sql

select "Booster_Version" as max_payload_mass_booster_versions
from SPACEXTABLE
where "PAYLOAD_MASS_KG_" in
(select max("PAYLOAD_MASS_KG_")
from SPACEXTABLE)
```

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

```
15]: max_payload_mass_booster_versions
```

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

- Listing the names of the booster\_versions which have carried the maximum payload mass

# 2015 Launch Records

## 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: SQLite 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.**

```
.6]: %%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'
```

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

```
.6]:
```

	Month	Landing_Outcome	Booster_Version	Launch_Site
	10	Failure (drone ship)	F9 v1.1 B1012	CCAFS LC-40
	04	Failure (drone ship)	F9 v1.1 B1015	CCAFS LC-40

- listing the records which will display the month names, failure landing\_outcomes in drone ship ,booster versions, launch site for the months in year 2015

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

## 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 Date, "Landing_Outcome", count("Landing_Outcome") as Landing_Outcome_count
from SPACEXTABLE
where Date between '2010-06-04' and '2017-03-20'
group by "Landing_Outcome"
order by Landing_Outcome_count desc
```

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

Done.

Date	Landing_Outcome	Landing_Outcome_count
2012-05-22	No attempt	10
2015-12-22	Success (ground pad)	5
2016-08-04	Success (drone ship)	5
2015-10-01	Failure (drone ship)	5
2014-04-18	Controlled (ocean)	3
2013-09-29	Uncontrolled (ocean)	2
2015-06-28	Precluded (drone ship)	1
2010-08-12	Failure (parachute)	1

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

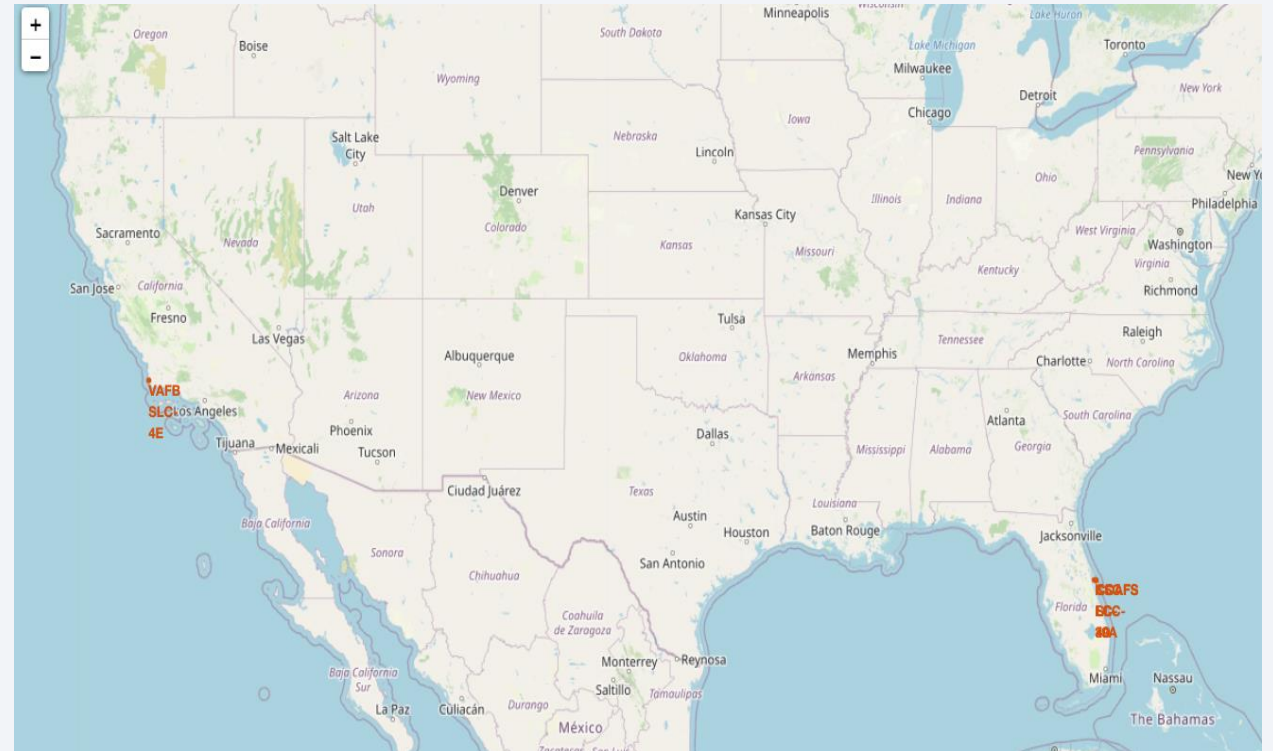
A satellite view of Earth from space, showing the curvature of the planet and city lights at night. The background is a deep blue gradient.

Section 3

# Launch Sites Proximities Analysis

# All Launches site location markers

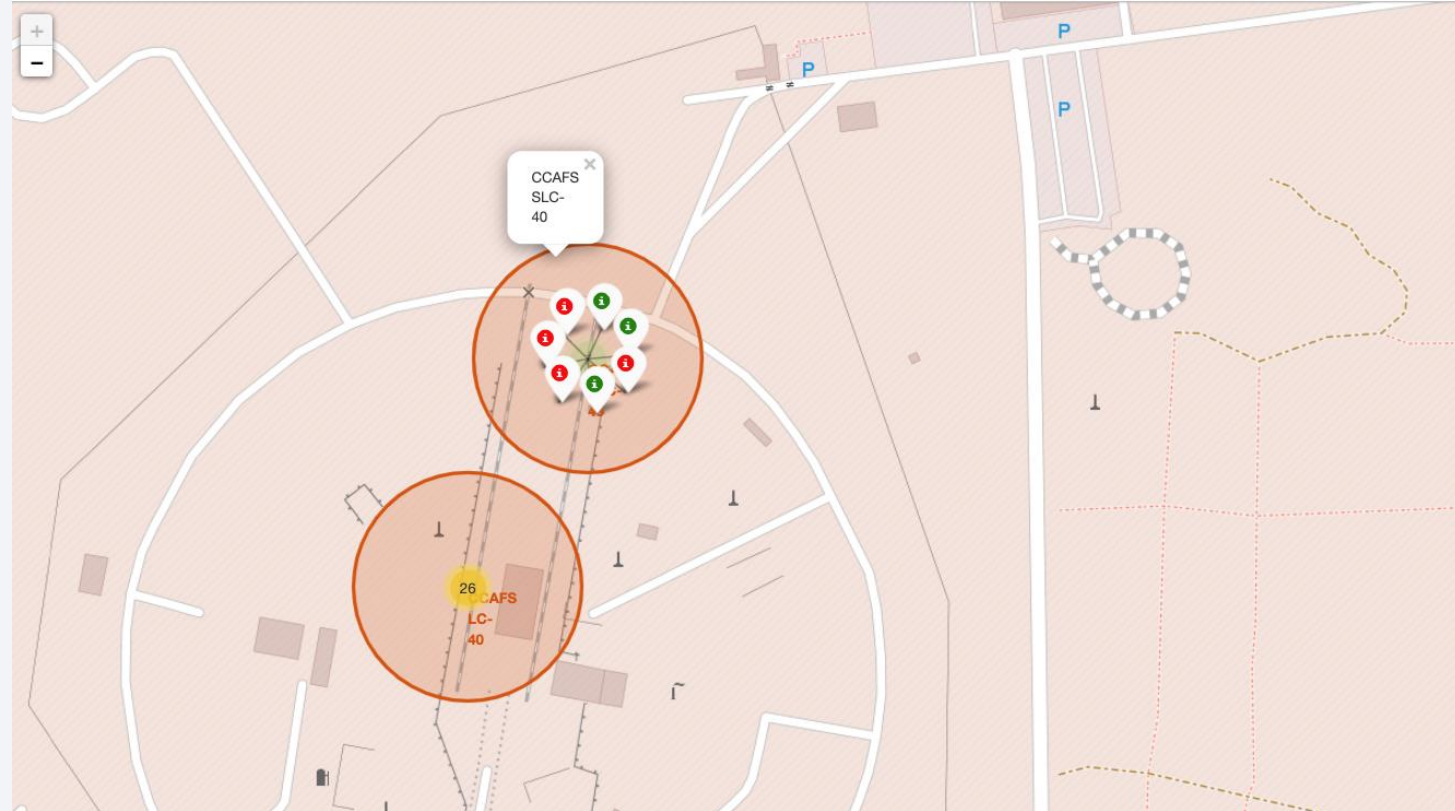
- All launch sites are in proximity with the equator line. The land is moving faster at the equator than any other place at the surface of the Earth. Anything on the surface of the Earth at the equator is already moving at 1670 km/hr. If a ship is launched from the equator it goes up into the space and it is also moving around the Earth at the same time 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 minimizes the risk of having any debris dropping or near people.





# Color-labeled launch records on the map

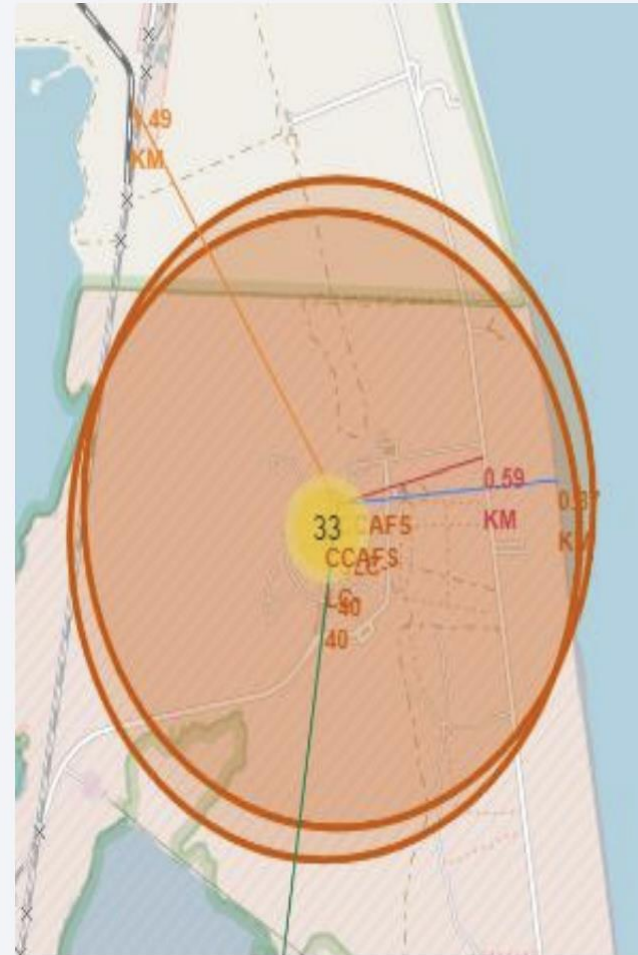
- Colored markers are used to help identify the launch sites with the highest success rate  
Green = Successful  
Red = Failed
- CCAFS SLC-49 has low success rate





## <Dashboard Screenshot 3>

- It is very close to railway (Nasa Railway: 1.49 km)
- It is very close to highway (Samuel C Phillips Highway: 1.49 km)
- It is very close to coastline (0.87 km)
- It is relatively close to its closest city (Cape Canaveral: 18.16 km)

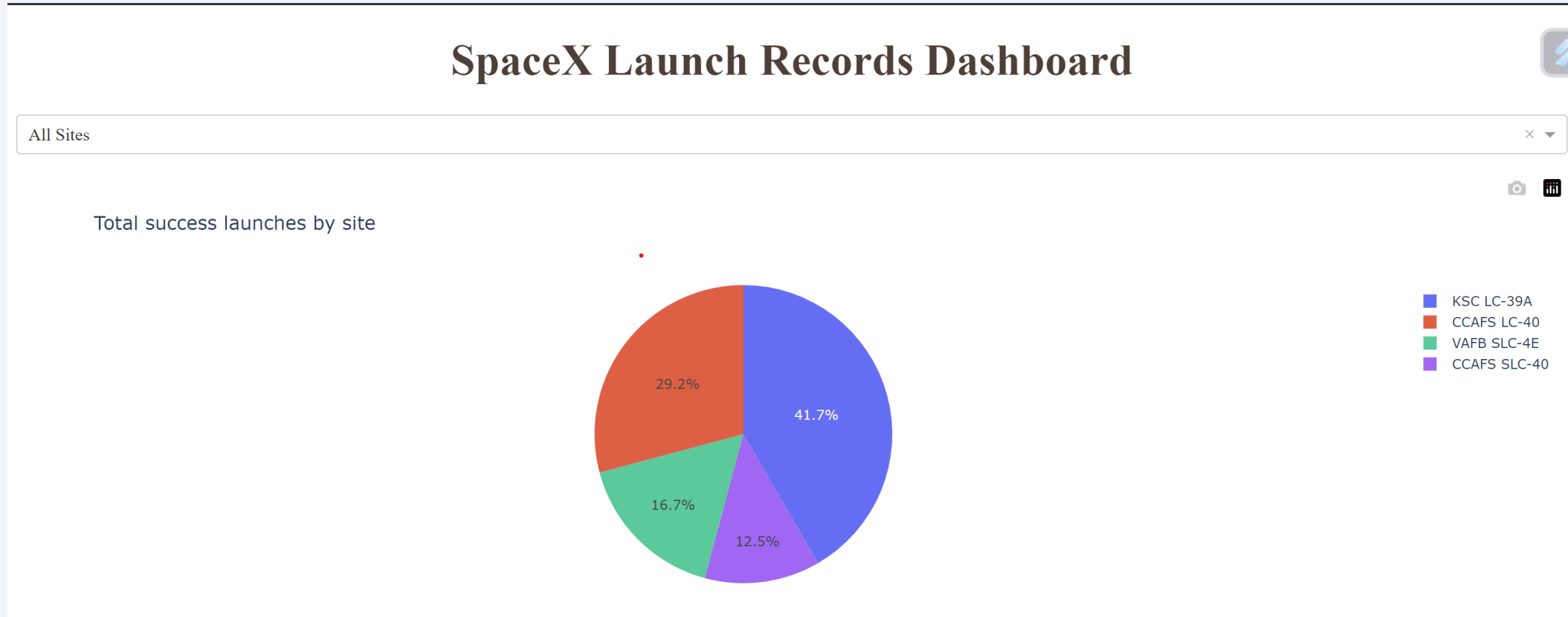




Section 4

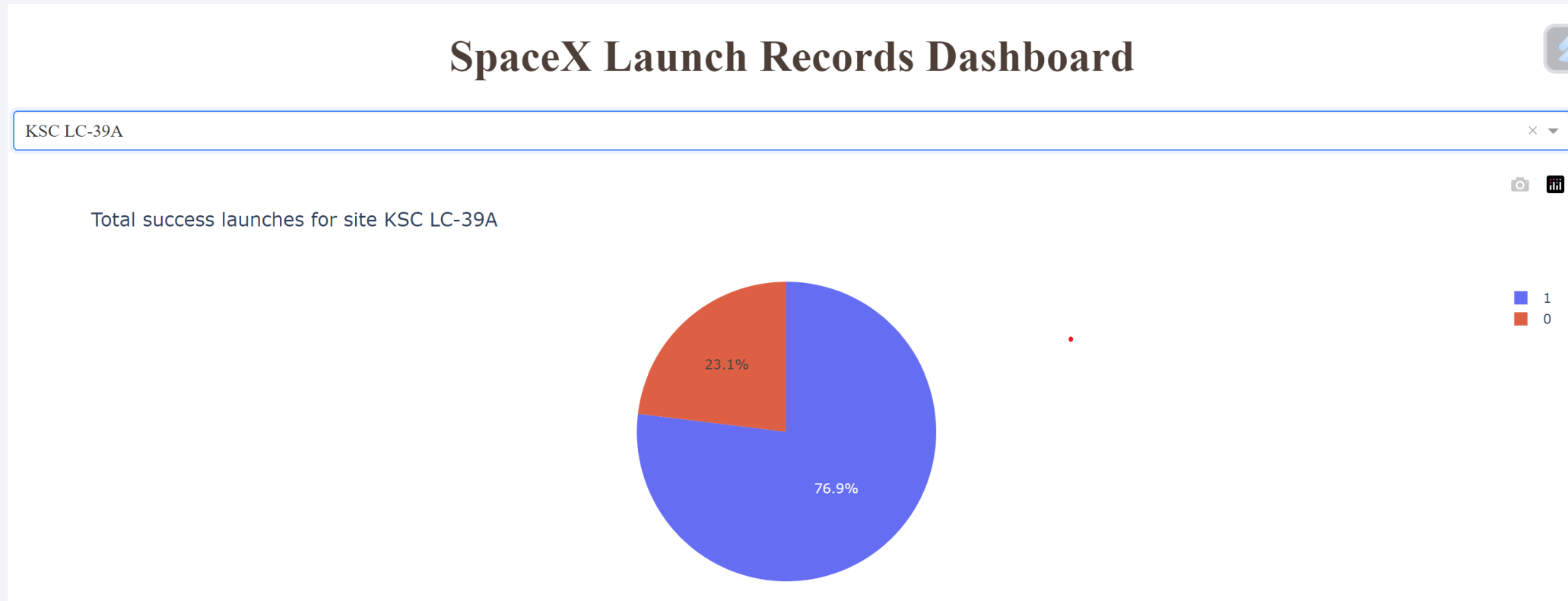
# Build a Dashboard with Plotly Dash

# Launch Success count for all sites



From the pie chart we can clearly see that that launch site KSC LC-39A has the highest success rate in launches.

# Success Launches for KSC LC-39A



KCL LC-39A has the highest success rate : 76.9%

# Payload Mass vs Launch Outcome for all sites

- Payloads between 2000 and 5000 kg have the highest success rate.
- Payloads between 5000 and 10000 kg have the lowest success rate.
- It is surprising that only one launch took place with the F9 Booster Version of type B5, which was successful. Therefore, Booster Version B5 has the highest launch success rate(100%)



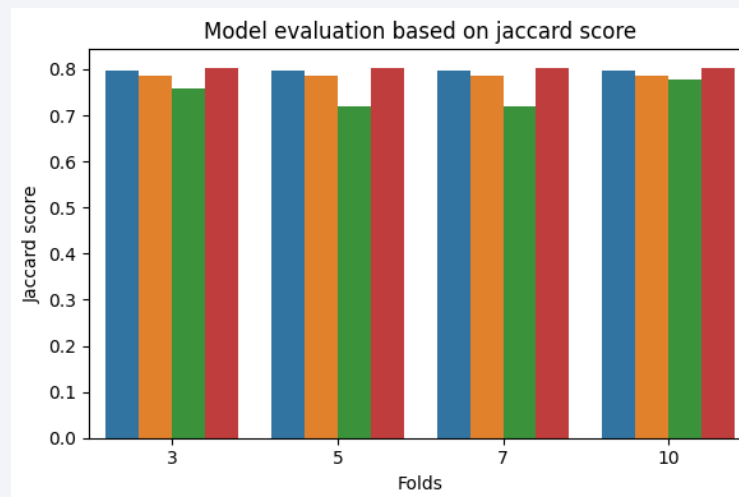
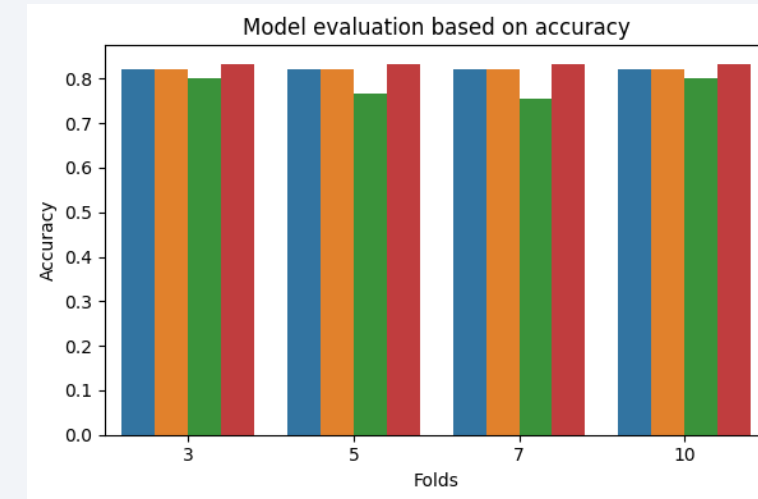
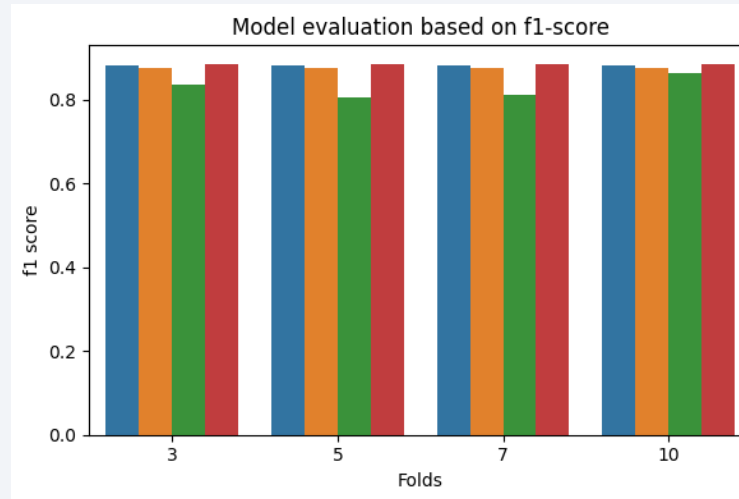


Section 5

# Predictive Analysis (Classification)

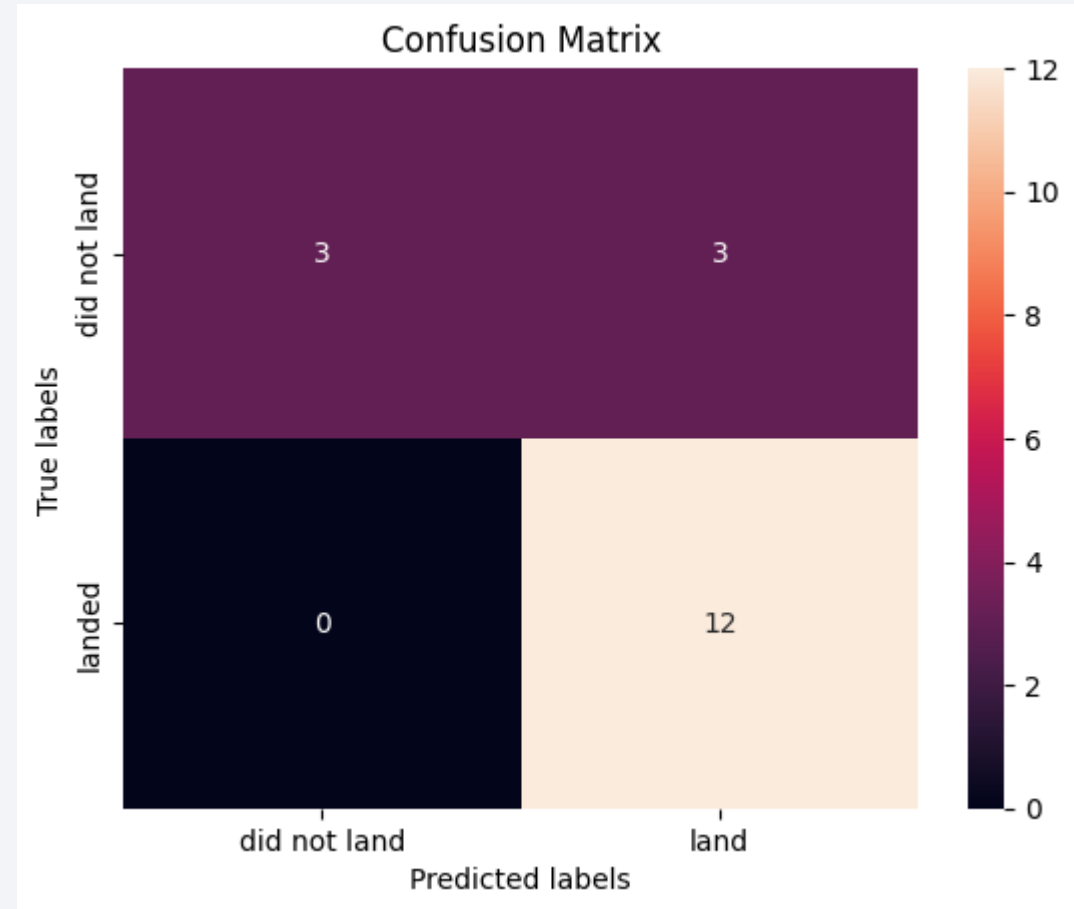
# Classification Accuracy

Because of the small test sample size we use K -Fold cross validation to determine the model that performs the best. From the bar plots it is obvious that KNN has the best performance across all folds and metrics .



# Confusion Matrix

Examining the confusion matrix. KNN can distinguish between the different classes. The major problem is false positives





# Conclusions

---

- K Nearest Neighbors is the best classification algorithm for this dataset.
- Launches with a low payload mass show better results than launches with a higher pK Nearest Neighbors is the best classification algorithm for this dataset.
- Launches with a low payload mass show better results than launches with a higher payload mass.
- Most of launch sites are in proximity to the Equator line and all launch sites are in very close proximity to the coastline.
- The launch success rate increases over the years.
- Orbits ES-L1,GEO,HEO and SSO have 100% success rate.
- Booster Version B5 has the highest success rate.ayload mass.

# Appendix

---

Thankyou to

- Coursea
- IBM
- etc

Thank you!

