

# Winning Space Race with Data Science

Bladimir Garcia Rosario  
09 July 2022



# Outline

---

- Executive Summary
- Introduction
- Methodology
- Results
- Conclusion

# Executive Summary

---

## Summary of methodologies:

- Data collection
- Data wrangling
- Exploratory Data Analysis with Data Visualization
- Exploratory Data Analyst 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 in screenshots
- Predictive analysis results

# Introduction

---

## Project background and context

SpaceX stands out as the leading player in the era of commercial space exploration, revolutionizing the affordability of space travel. On their website, they proudly showcase Falcon 9 rocket launches at a remarkably lower price tag of 62 million dollars compared to other providers who charge a hefty 165 million dollars per launch. A key factor contributing to this cost advantage is SpaceX's ability to reuse the first stage of their rockets. By accurately predicting the landing success of the first stage, we can reliably estimate the overall cost of a launch. Leveraging both publicly available information and advanced machine learning models, our goal is to forecast whether SpaceX will successfully recover and reuse the first stage.

## Problems you want to find answers

- In what ways do factors like payload mass, launch site, number of flights, and orbits impact the likelihood of a successful first stage landing?
- Is there a noticeable upward trend in the rate of successful landings over the years?
- Which algorithm would be most suitable for binary classification in this particular scenario?

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

---

The data collection process consisted of gathering information from two primary sources: the SpaceX REST API and a table in SpaceX's Wikipedia entry. Both of these data collection methods were utilized to ensure comprehensive data retrieval for a thorough analysis of the launches.

## **Data Columns are obtained by using SpaceX REST API:**

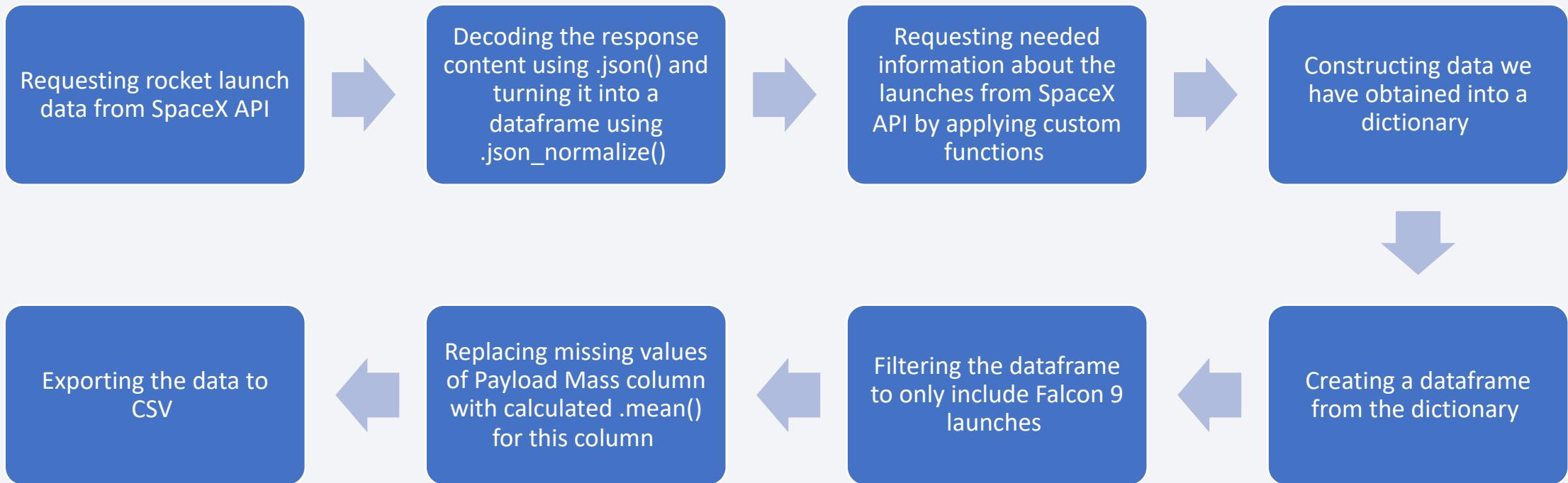
- FlightNumber
- Date
- BoosterVersion
- PayloadMass
- Orbit
- LaunchSite
- Outcome
- Flights
- GridFins
- Reused
- Legs
- LandingPad
- Block
- ReusedCount
- Serial
- Longitude
- Latitude

## **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

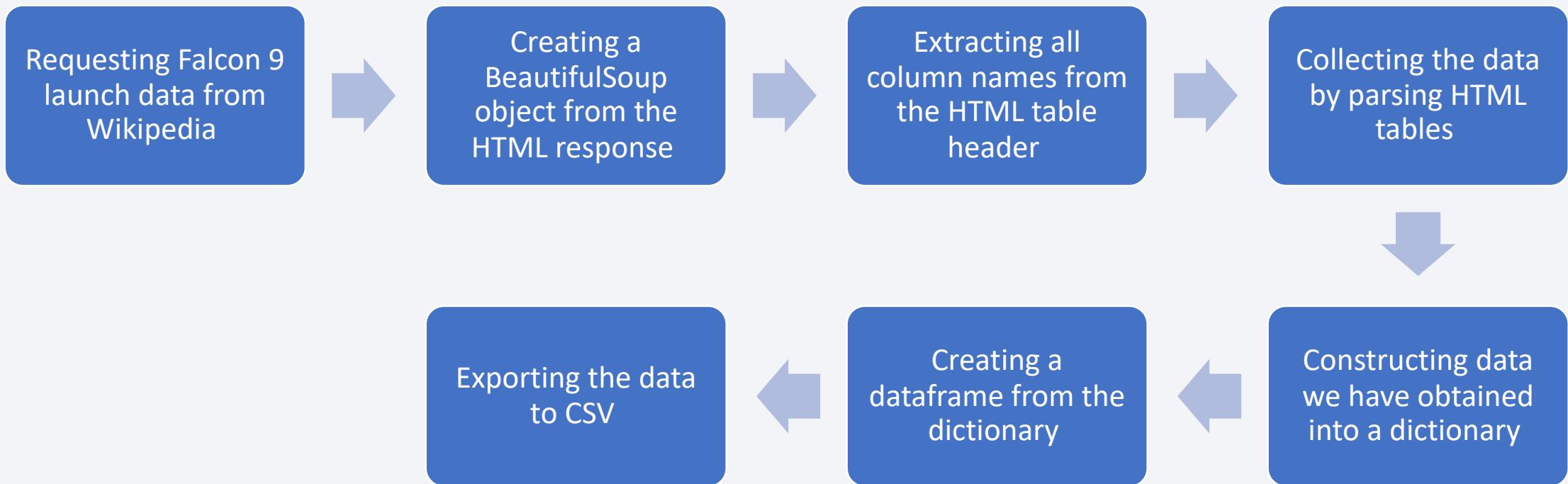
# Data Collection – SpaceX API

---



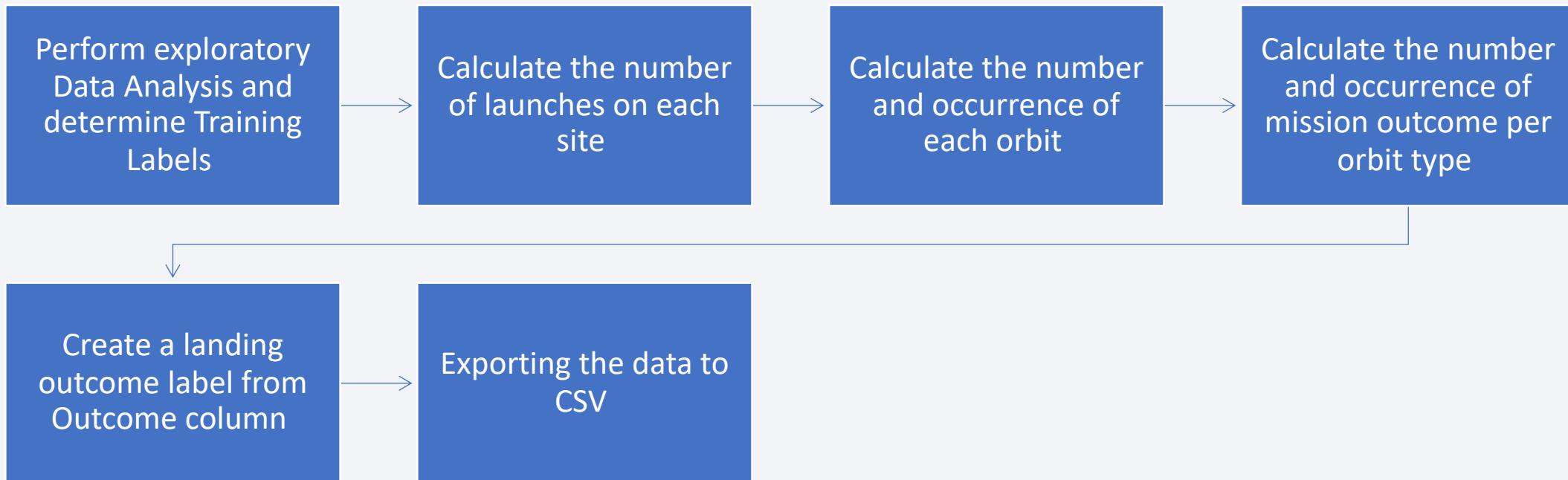
# Data Collection - Scraping

---



# Data Wrangling

Within the dataset, there are various scenarios in which the landing of the booster was not successful. In some cases, landing attempts were made but failed due to unforeseen circumstances. For instance, when the mission outcome is labeled as "True Ocean," it signifies that the landing was successfully executed in a specific area of the ocean. On the other hand, a label of "False Ocean" indicates an unsuccessful landing in the designated ocean region. Similarly, "True RTLS" signifies a successful landing on a ground pad, while "False RTLS" represents an unsuccessful ground pad landing. In the case of "True ASDS," it denotes a successful landing on a drone ship, whereas "False ASDS" indicates an unsuccessful landing on a drone ship. These outcomes have been transformed into training labels, where a value of "1" indicates a successful booster landing and a value of "0" indicates an unsuccessful landing.



# EDA with Data Visualization

---

## Chart 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 offer a visual representation that helps us understand how variables relate to each other. This understanding is valuable when integrating them into machine learning models.

When it comes to comparing distinct categories, bar charts provide a powerful tool. They allow us to explore the connection between specific categories and the corresponding measurements, leading to meaningful insights.

For tracking data over time and uncovering patterns or changes, line charts prove to be highly effective. They enable us to visualize trends in time series data, facilitating analysis and providing valuable insights into the underlying patterns.

# 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:

Implemented a marker on the map, incorporating a circle, a popup label, and a text label to represent the NASA Johnson Space Center. This marker was positioned based on the center's latitude and longitude coordinates, serving as the starting location for further exploration.

Integrated markers on the map for all launch sites, each accompanied by a circle, a popup label, and a text label. These markers were placed according to the latitude and longitude coordinates of the respective launch sites, effectively visualizing their geographical positions and their proximity to the Equator and coastlines.

- Markers of all Launch Sites:

Utilized a Marker Cluster feature to display colored markers representing successful launches (green) and failed launches (red). This visual representation allowed for easy identification of launch sites with comparatively higher success rates, based on the clustering of green markers.

- Markers of all Launch Sites:

Incorporated colored lines on the map to visually represent the distances between the launch site KSC LC-39A and its surrounding areas such as the railway, highway, coastline, and closest city. This addition allows for a clear visualization of the proximity of the launch site to these features.

# Build a Dashboard with Plotly Dash

---

## **Launch Sites Dropdown List:**

Implemented a dropdown menu feature that allows for the selection of a specific launch site. This interactive functionality enhances the user experience by providing the ability to easily choose and focus on a particular launch site of interest.

## **Launch Sites Dropdown List:**

Incorporated a pie chart into the visualization to display the overall count of successful launches across all launch sites. Additionally, when a specific launch site is selected, the pie chart provides a breakdown of the success and failure counts for that particular site. This visual representation enhances the understanding of the success rates at individual launch sites and allows for easy comparison between sites.

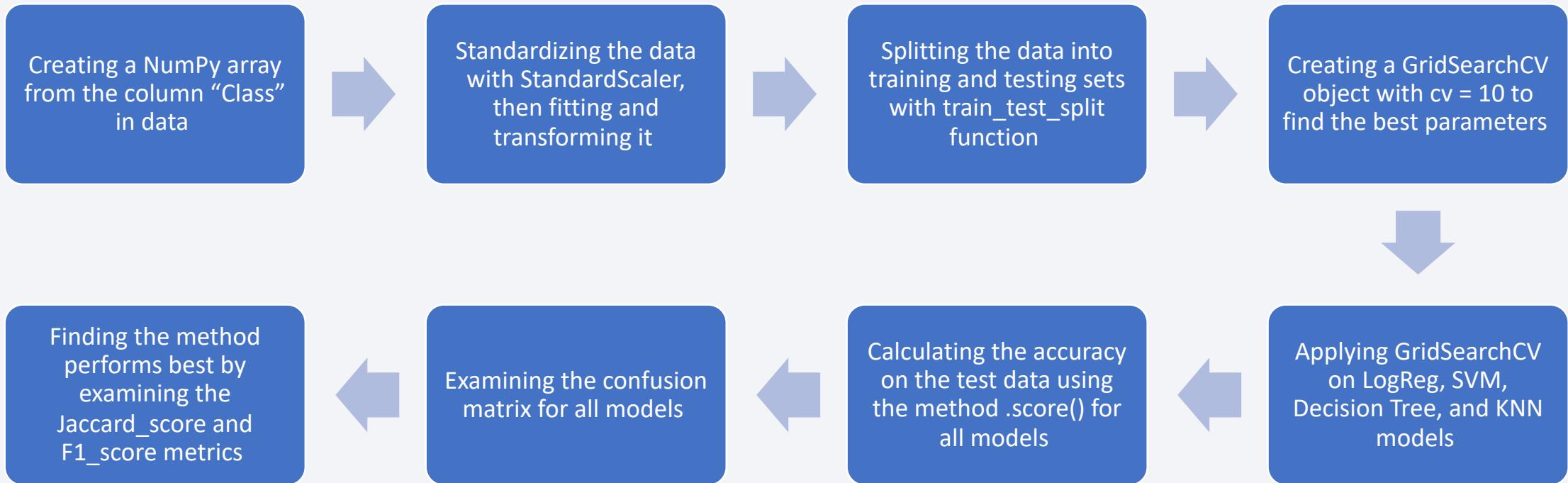
## **Launch Sites Dropdown List:**

Added a slider to select Payload range.

## **Launch Sites Dropdown List:**

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

# Predictive Analysis (Classification)



# Results

---

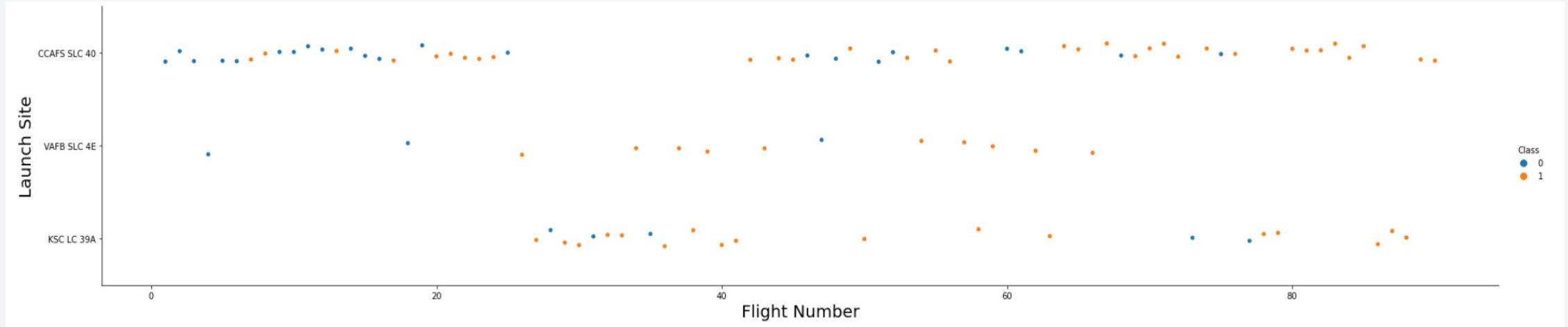
- Exploratory data analysis results
- Interactive analytics demo in screenshots
- Predictive analysis results

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

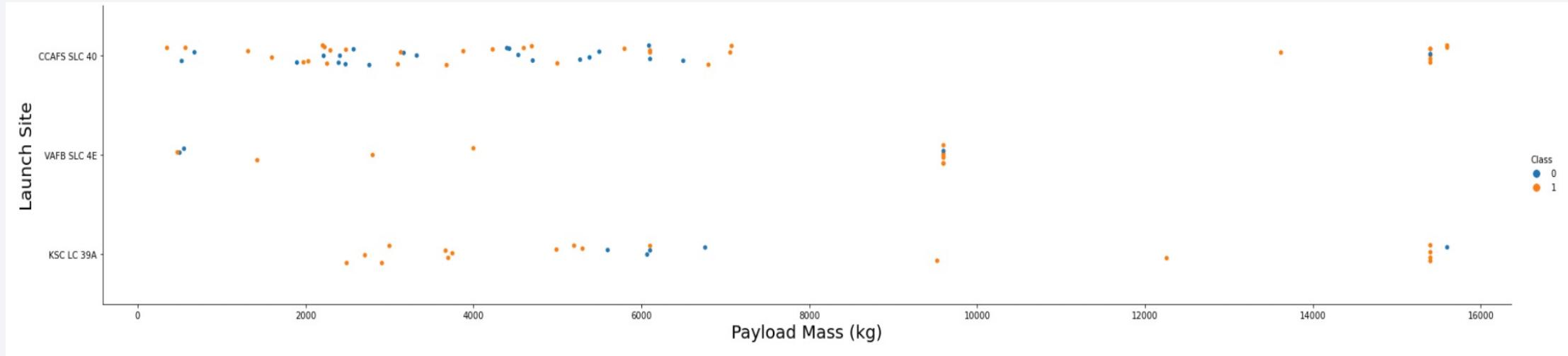
# Flight Number vs. Launch Site



## 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

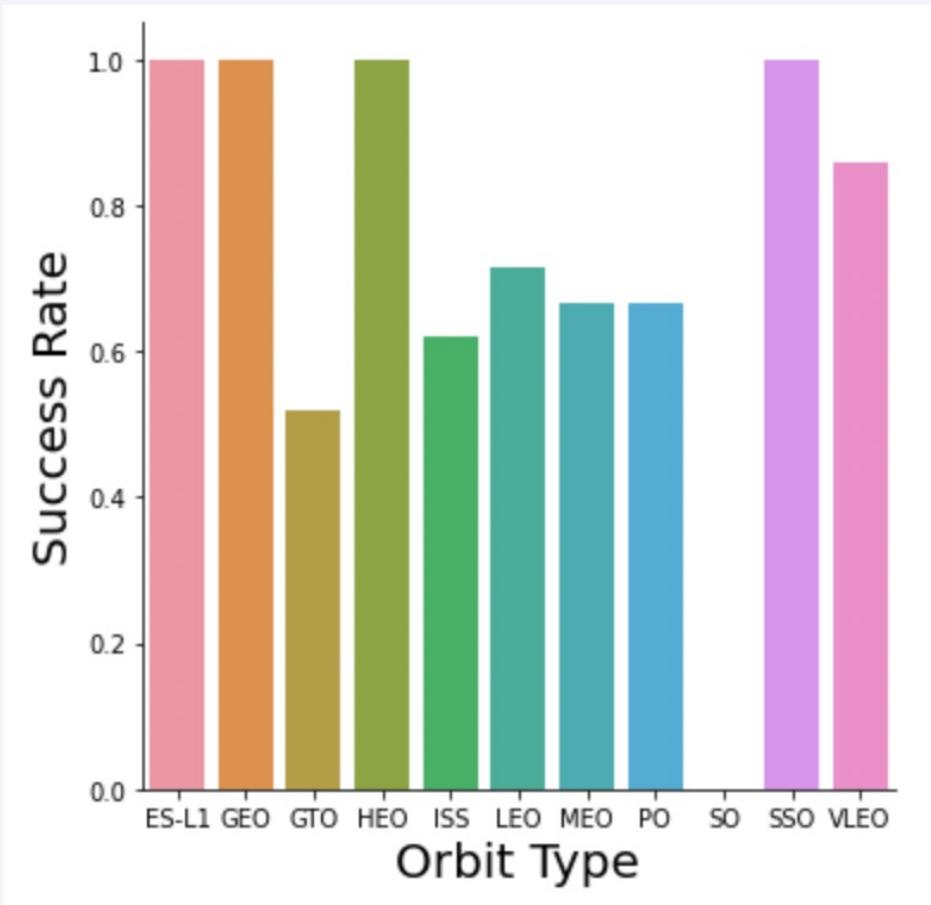


## Explanation:

- 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:

### Orbits with 100% success rate:

- ES-L1, GEO, HEO, SSO

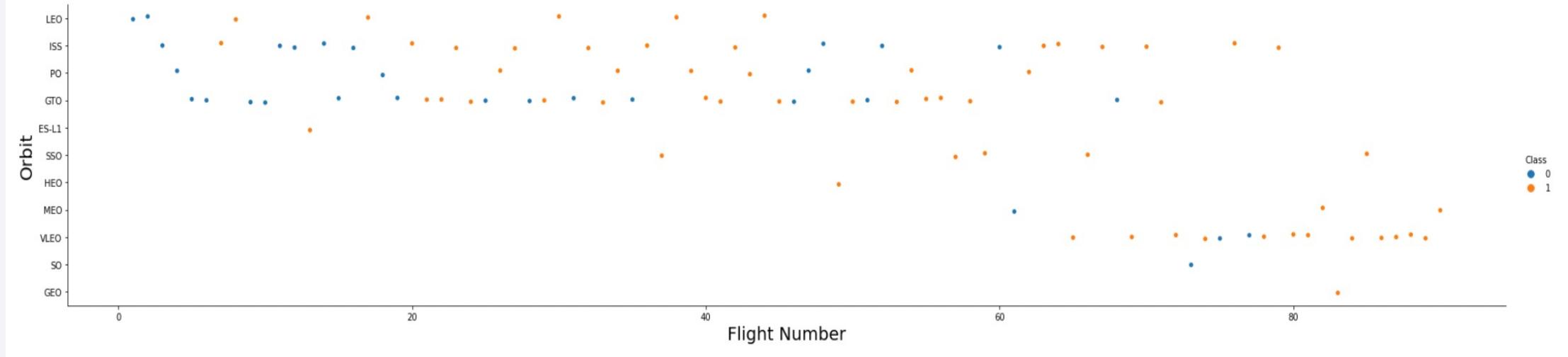
### Orbits with 0% success rate:

- SO

### Orbits with success rate between 50% and 85%:

- GTO, ISS, LEO, MEO, PO

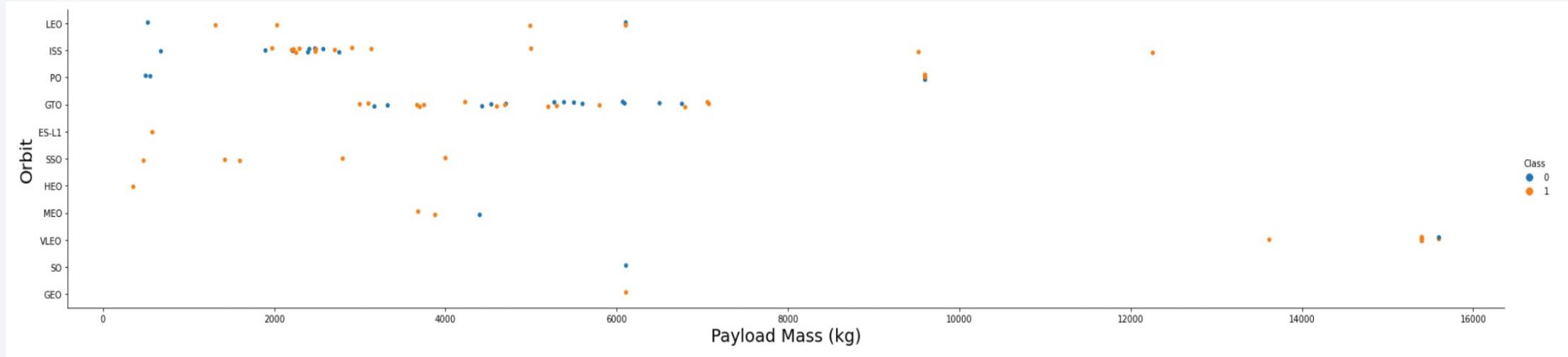
# Flight Number vs. Orbit Type



## 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

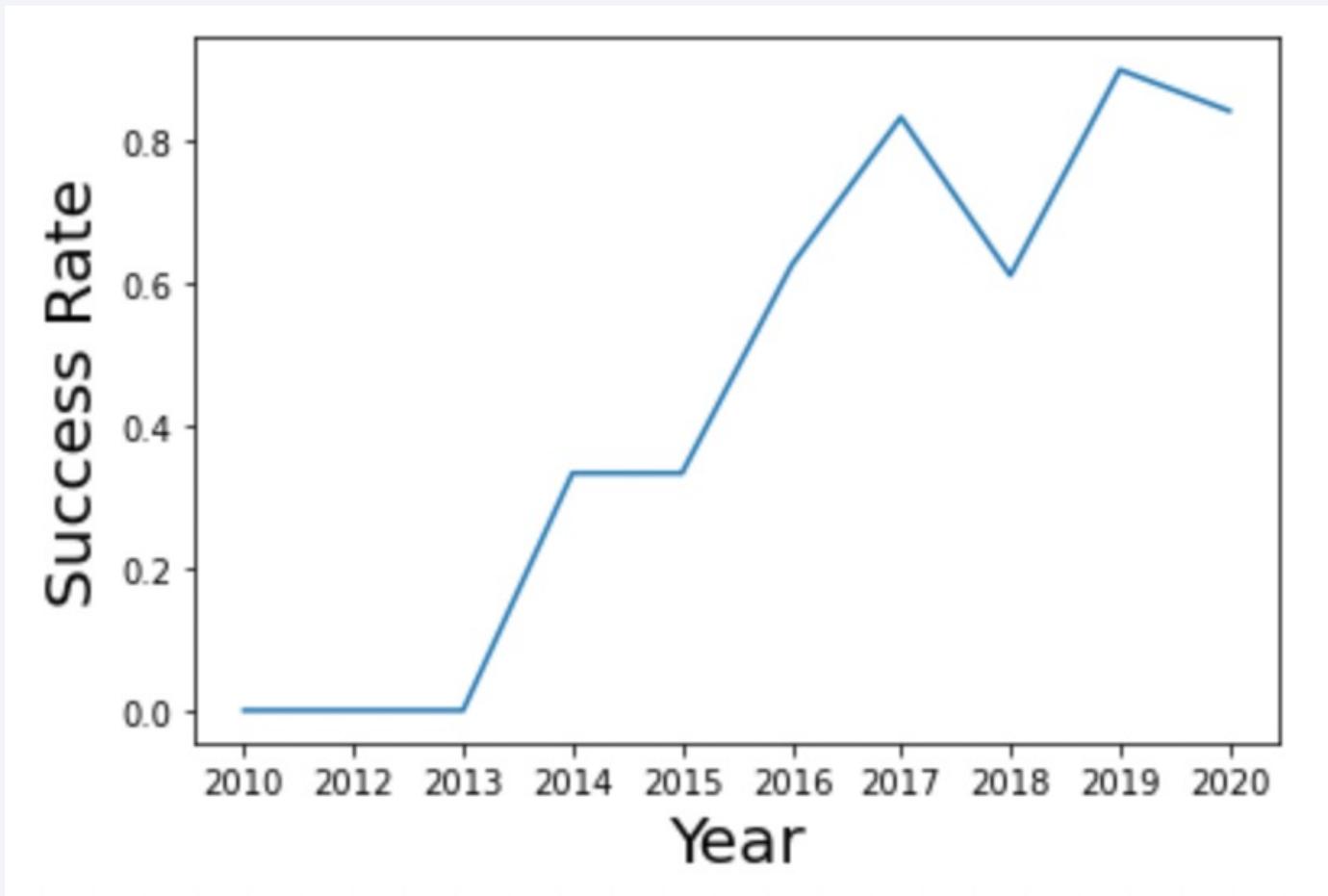


## Explanation:

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.

# All Launch Site Names

---

```
In [4]: %sql select distinct launch_site from SPACEXDATASET;
```

```
* ibm_db_sa://wzf08322:***@0c77d6f2-5da9-48a9-81f8-86b520b87518.bs2io90108kqb1od8lcg.databases.appdomain.cloud:31198/bludb  
Done.
```

Out[4]:

launch_site
CCAFS LC-40
CCAFS SLC-40
KSC LC-39A
VAFB SLC-4E

## 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

## Explanation:

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

# Total Payload Mass

---

```
In [6]: %sql select sum(payload_mass_kg_) as total_payload_mass from SPACEXDATASET where customer = 'NASA (CRS)';

* ibm_db_sa://wzf08322:***@0c77d6f2-5da9-48a9-81f8-86b520b87518.bs2io90108kqb1od8lcg.databases.appdomain.cloud:31198/bludb
Done.
```

```
Out[6]: total_payload_mass
45596
```

## Explanation:

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

# Average Payload Mass by F9 v1.1

---

```
In [7]: %sql select avg(payload_mass_kg_) as average_payload_mass from SPACEXDATASET where booster_version like '%F9 v1.1%';  
* ibm_db_sa://wzf08322:***@0c77d6f2-5da9-48a9-81f8-86b520b87518.bs2io90108kqb1od81cg.databases.appdomain.cloud:31198/bludb  
Done.
```

```
Out[7]:
```

average_payload_mass
2534

## Explanation:

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

# First Successful Ground Landing Date

---

```
In [8]: %sql select min(date) as first_successful_landing from SPACEXDATASET where landing__outcome = 'Success (ground pad)';  
* ibm_db_sa://wzf08322:***@0c77d6f2-5da9-48a9-81f8-86b520b87518.bs2io90108kqb1od81cg.databases.appdomain.cloud:31198/bludb  
Done.  
Out[8]: first_successful_landing  
2015-12-22
```

## Explanation:

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

# Successful Drone Ship Landing with Payload between 4000 and 6000

---

```
In [9]: %sql select booster_version from SPACEXDATASET where landing__outcome = 'Success (drone ship)' and payload_mass__kg_ between 4000 and 6000;  
* ibm_db_sa://wzf08322:***@0c77d6f2-5da9-48a9-81f8-86b520b87518.bs2io90108kqb1od81cg.databases.appdomain.cloud:31198/bludb  
Done.
```

Out[9]:

booster_version
F9 FT B1022
F9 FT B1026
F9 FT B1021.2
F9 FT B1031.2

## Explanation:

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

# Total Number of Successful and Failure Mission Outcomes

---

```
In [10]: %sql select mission_outcome, count(*) as total_number from SPACEXDATASET group by mission_outcome;  
* ibm_db_sa://wzf08322:***@0c77d6f2-5da9-48a9-81f8-86b520b87518.bs2io90108kqb1od8lcg.databases.appdomain.cloud:31198/bludb  
Done.
```

Out[10]:

mission_outcome	total_number
Failure (in flight)	1
Success	99
Success (payload status unclear)	1

## Explanation:

Listing the total number of successful and failure mission outcomes

# Boosters Carried Maximum Payload

```
In [11]: %sql select booster_version from SPACEXDATASET where payload_mass_kg_ = (select max(payload_mass_kg_) from SPACEXDATASET);  
* ibm_db_sa://wzf08322:***@0c77d6f2-5da9-48a9-81f8-86b520b87518.bs2io90108kqb1od81cg.databases.appdomain.cloud:31198/bludb  
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

## Explanation:

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

# 2015 Launch Records

```
In [12]: %%sql select monthname(date) as month, date, booster_version, launch_site, landing__outcome from SPACEXDATASET  
where landing__outcome = 'Failure (drone ship)' and year(date)=2015;  
* ibm_db_sa://wzf08322:***@0c77d6f2-5da9-48a9-81f8-86b520b87518.bs2io90108kqb1od81cg.databases.appdomain.cloud:31198/bludb  
Done.
```

Out[12]:

MONTH	DATE	booster_version	launch_site	landing__outcome
January	2015-01-10	F9 v1.1 B1012	CCAFS LC-40	Failure (drone ship)
April	2015-04-14	F9 v1.1 B1015	CCAFS LC-40	Failure (drone ship)

## Explanation:

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

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

```
In [13]: %%sql select landing_outcome, count(*) as count_outcomes from SPACEXDATASET  
        where date between '2010-06-04' and '2017-03-20'  
        group by landing_outcome  
        order by count_outcomes desc;  
  
* ibm_db_sa://wzf08322:***@0c77d6f2-5da9-48a9-81f8-86b520b87518.bs2io90108kqb1od81cg.databases.appdomain.cloud:31198/bludb  
Done.
```

Out[13]:

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

## 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.

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

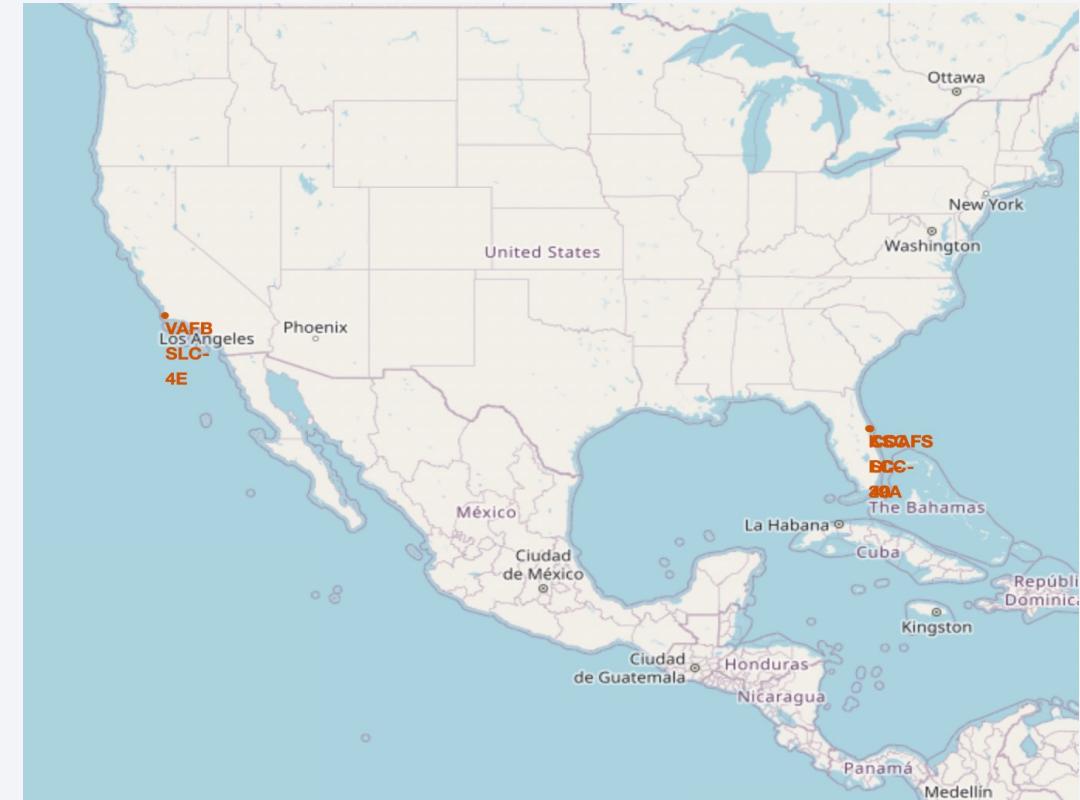
# All launch sites location markers on a global map

---

## 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.



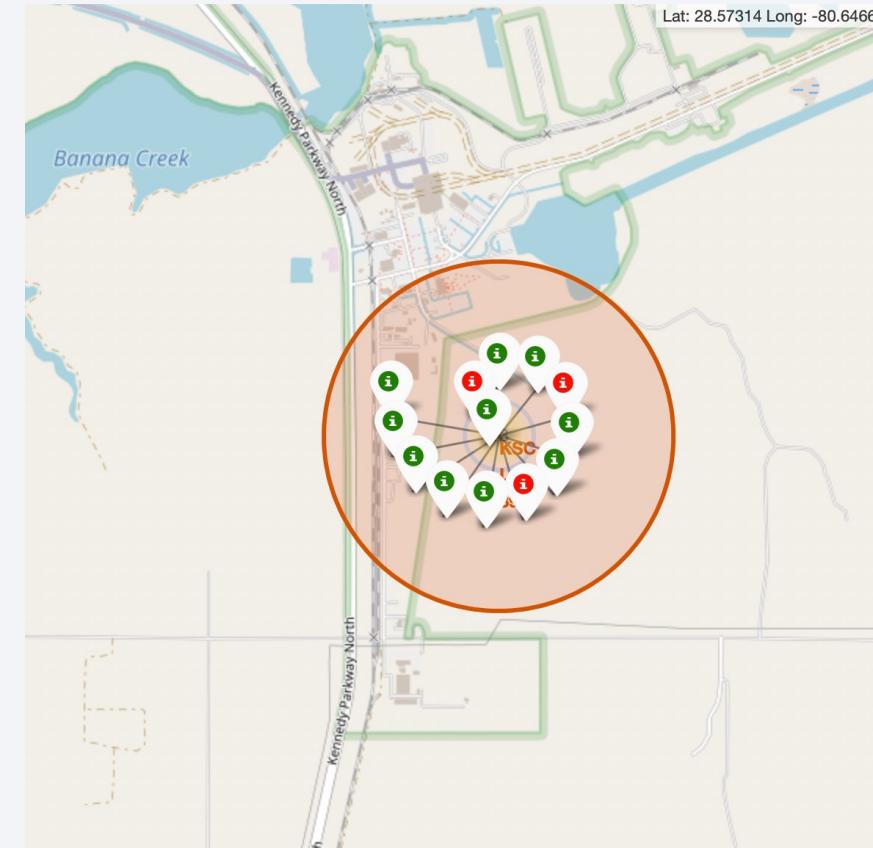
# Colour-labeled launch records on the map

## 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

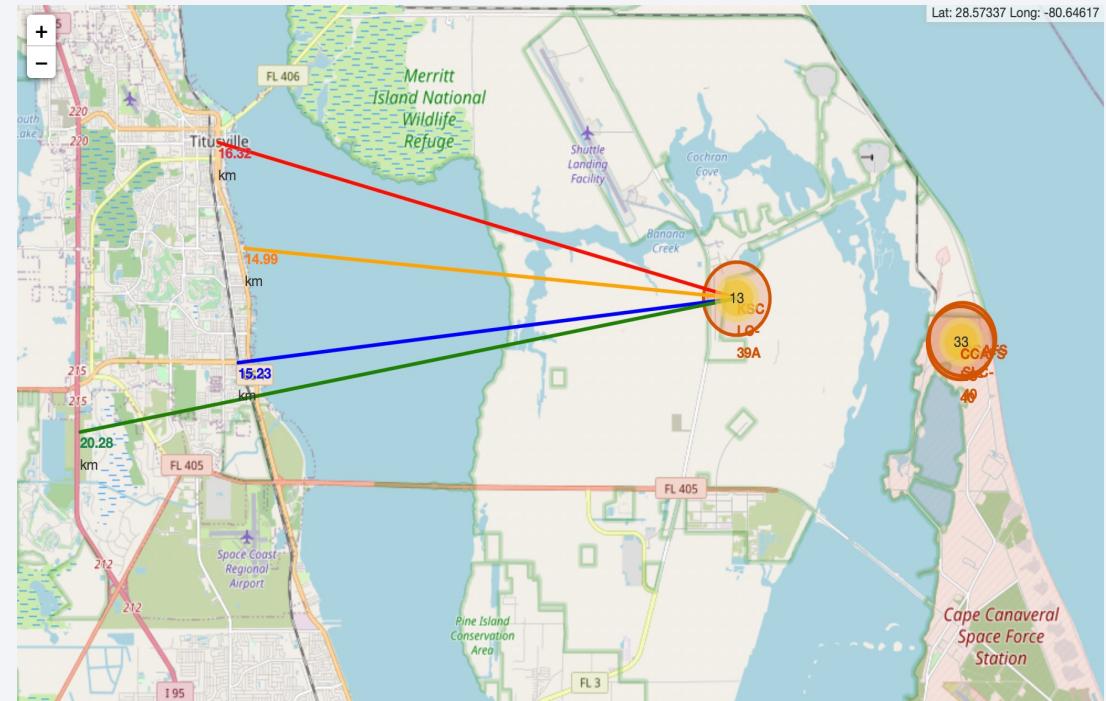
Launch Site KSC LC-39A has a very high Success Rate.



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

## 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.



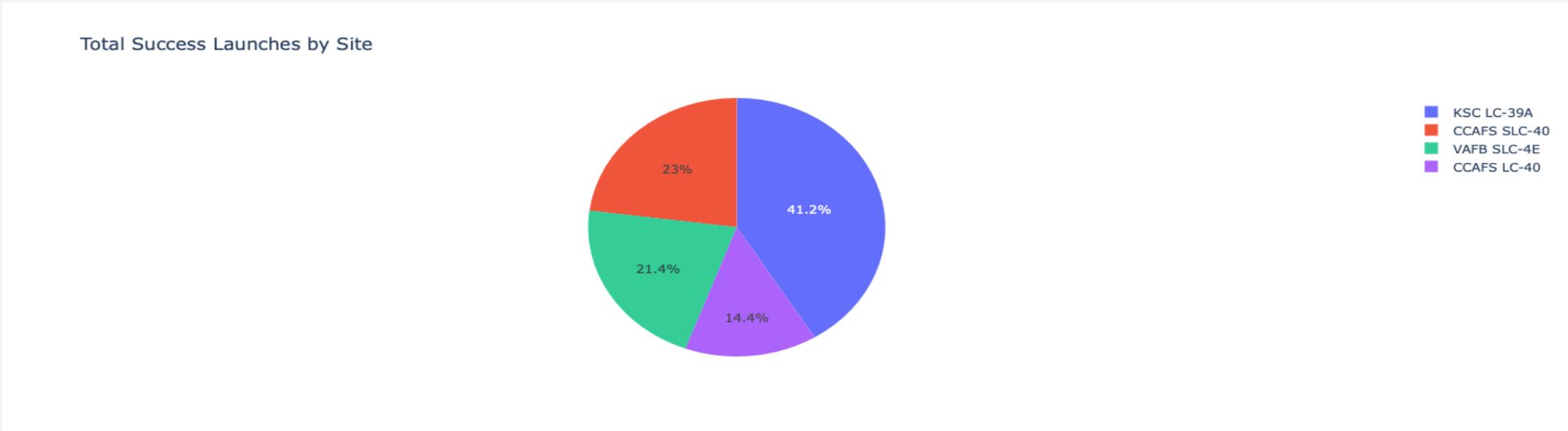
Section 4

# Build a Dashboard with Plotly Dash



# Launch success count for all sites

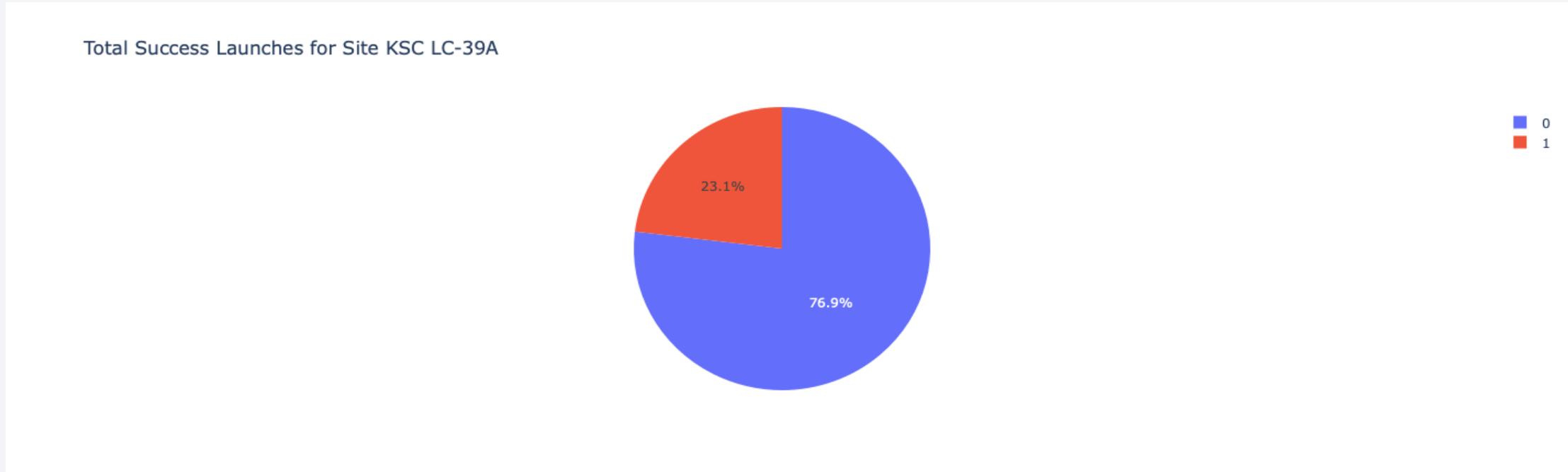
---



## 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



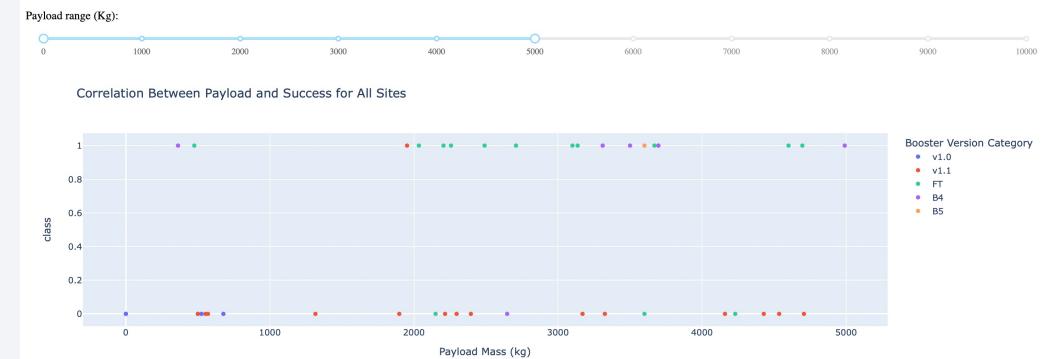
## Explanation:

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 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

---

## 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.

## Scores and Accuracy of the Test Set

	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

## Scores and Accuracy of the Entire Data Set

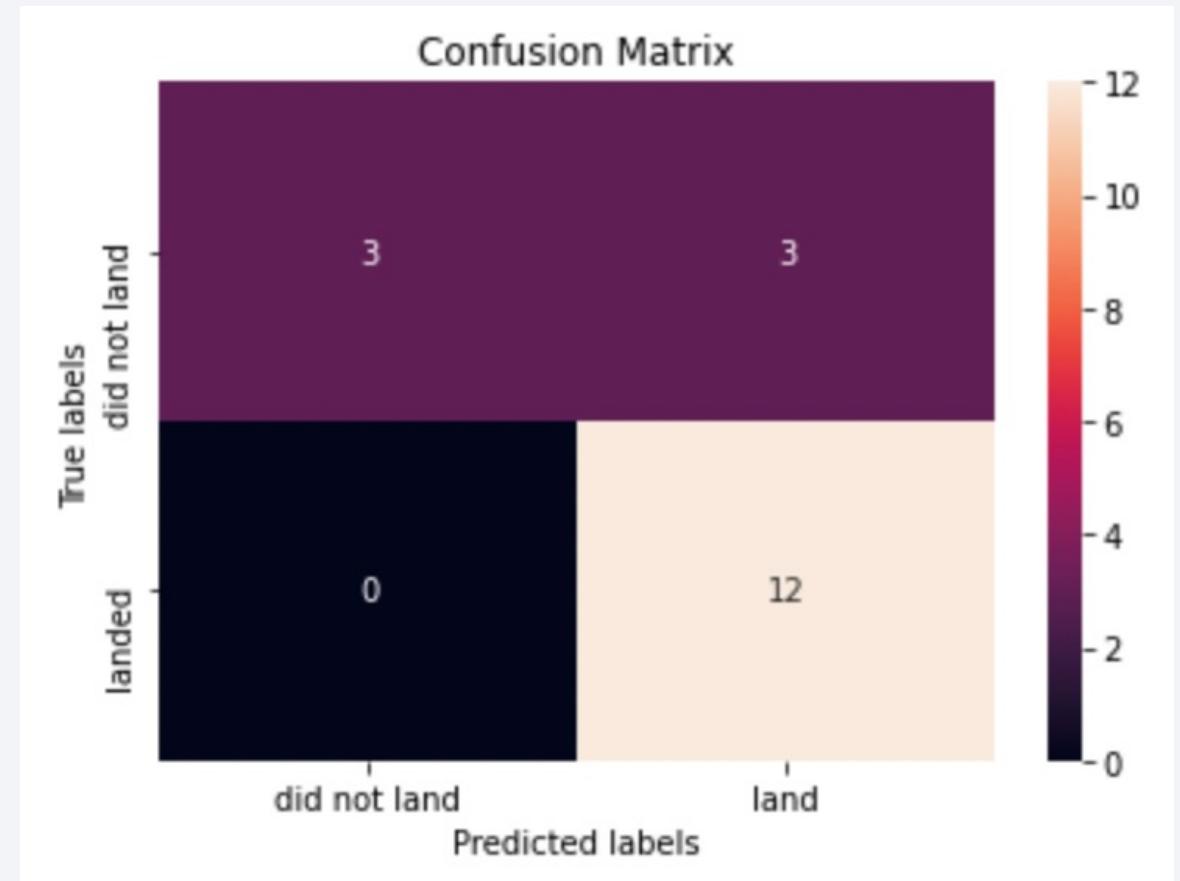
	LogReg	SVM	Tree	KNN
Jaccard_Score	0.833333	0.845070	0.882353	0.819444
F1_Score	0.909091	0.916031	0.937500	0.900763
Accuracy	0.866667	0.877778	0.911111	0.855556

# 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.

		Predicted Values	
		Negative	Positive
Actual Values	Negative	TN	FP
	Positive	FN	TP



# Conclusions

---

- 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.
- KSC LC-39A has the highest success rate of the launches from all the sites.
- Orbit ES-L1, GEO, HEO and SSO have 100% success rate.

Thank you!

