



IBM Developer  
SKILLS NETWORK

# Winning Space Race with Data Science

RYAN R. ALINCASTRE  
06.02.2025



# 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 in Screenshots
  - 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
  - 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?
  - What is the best algorithm that can be used for binary classification in this case?



Section 1

# Methodology

# Methodology

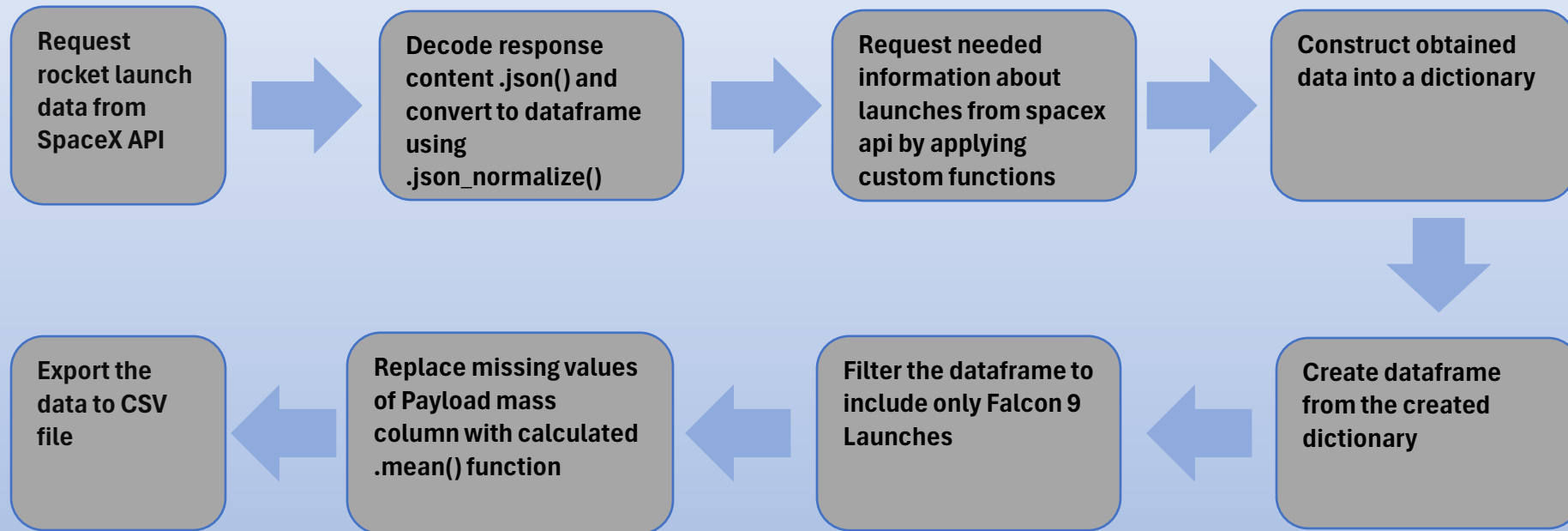
## Executive Summary

- Data collection methodology:
  - Using SpaceX Rest API
  - Using Web Scraping from Wikipedia
- Perform data wrangling
  - Filtering the data
  - Dealing with missing values
  - Using One Hot Encoding to prepare the
- 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 find best results

# Data Collection

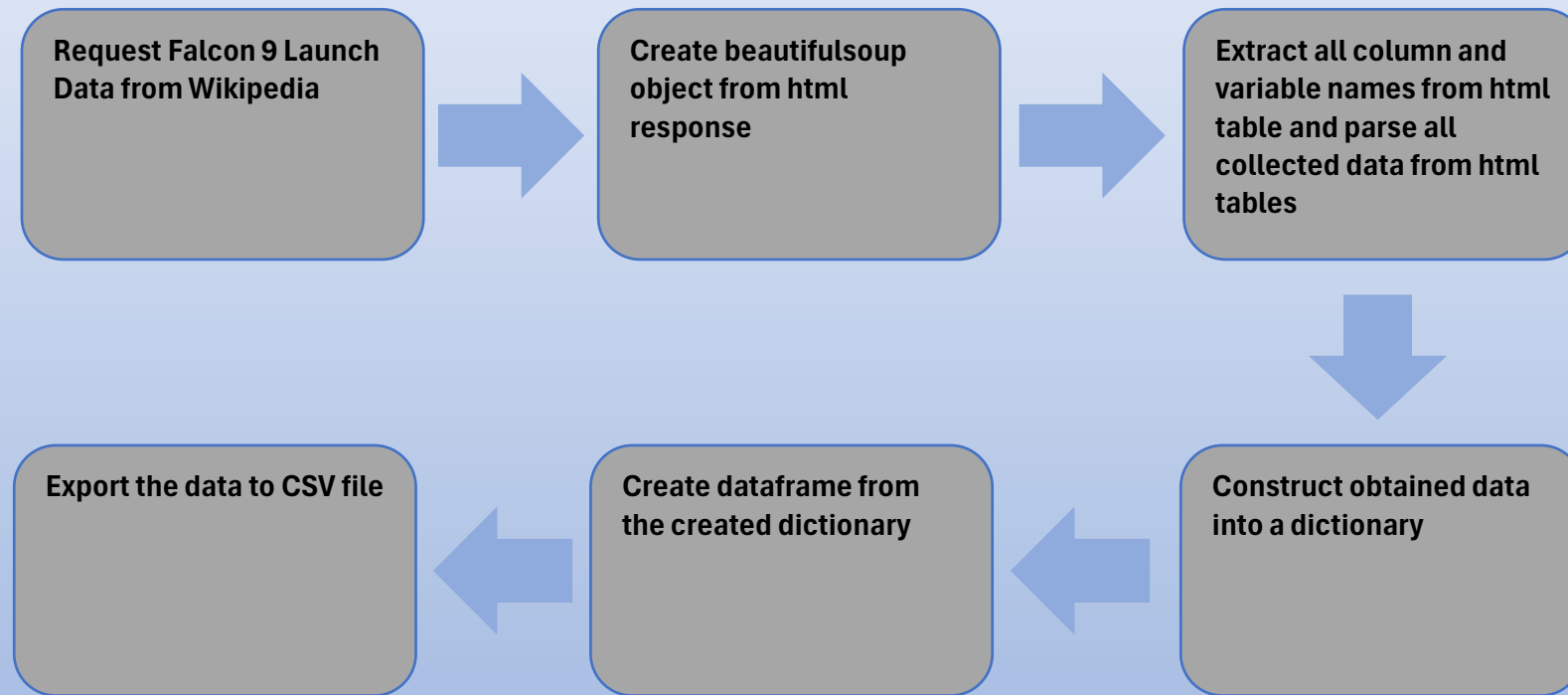
- Data collection process involved a combination of API requests from SpaceX REST API and Web Scraping data from a table in SpaceX's Wikipedia.
- REST API and Web Scraping are both used in data collection methods in order to get complete information about the launches for a more detailed analysis
- Data Columns 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

- 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.

Perform Exploratory Data Analysis and Determine Training Labels



Calculate the number of Launches on each site

Calculate the number and occurrence of each orbit

Calculate the number and occurrence of missing outcome per orbit type

Create a Landing Outcome label from Outcome Column

Export the data to CSV file

# EDA with Data Visualization

- Charts Plotted:
  - Scatter Plot for Flight Number vs. Payload Mass, Flight Number vs. Launch Site, Payload Mass vs. Launch Site, Flight Number vs. Orbit Type, and Payload Mass vs Orbit Type.
    - Scatter plots show the relationship between variables. If a relationship exists, they could be used in machine learning model.
  - Bar Chart for Success Rate vs. Orbit Type.
    - 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 Chart for Success Rate Yearly Trend.
    - Line Chart 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

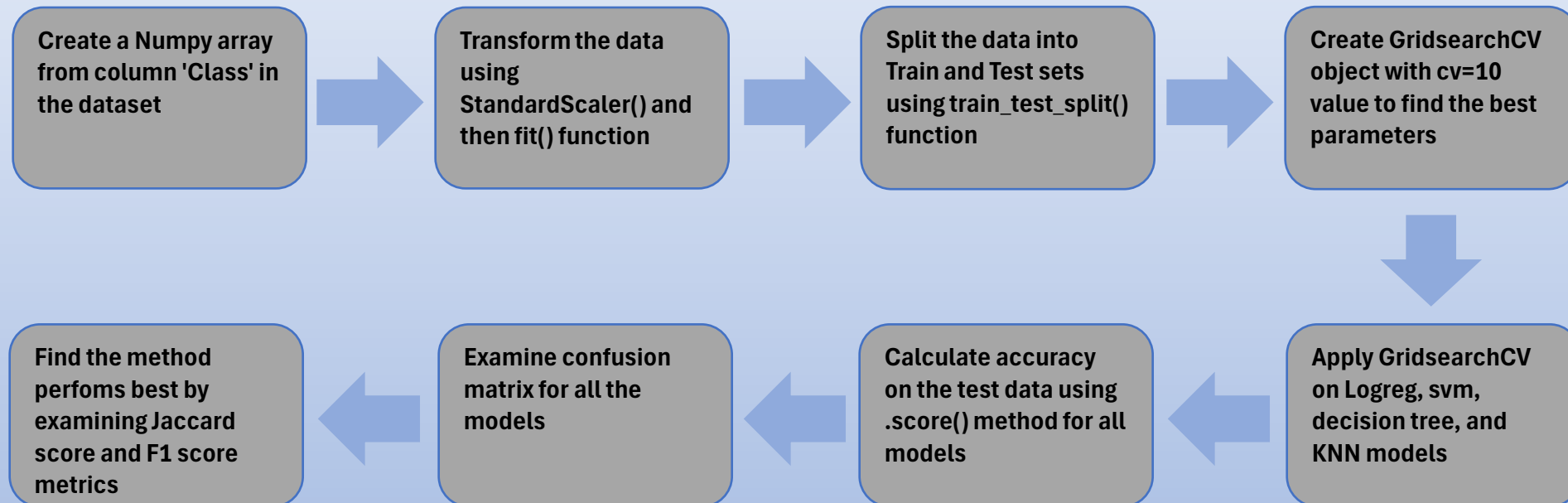
- Markers of all Launch Sites:
  - Added Marker with Circle, Popup Label and Text Label of NASA Johnson Space Center using its latitude and longitude coordinates as a start location.
  - Added Markers with Circle, Popup Label and Text Label of all Launch Sites using their latitude and longitude coordinates to show their geographical locations and proximity to Equator and coasts.
- Colored Markers of the launch outcomes for each Launch Site:
  - Added colored 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 colored Lines to show distances between the Launch Site KSC LC-39A (as an example) and its proximities like Railway, Highway, Coastline and Closest City.



# 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.

# Predictive Analysis (Classification)



# Results

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



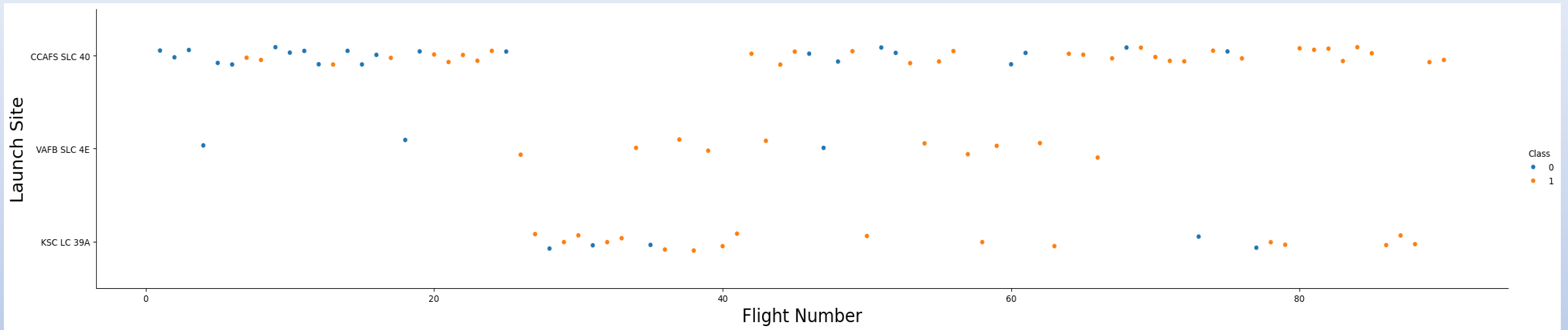
The background of the slide is an abstract composition. It features a dark blue base color. Overlaid on this are numerous diagonal streaks in shades of red and cyan, creating a sense of motion and depth. A faint, light blue grid pattern is also visible, particularly in the lower-left quadrant. The overall effect is high-tech and digital.

Section 2

# Insights drawn from EDA



# Flight Number vs. Launch Site

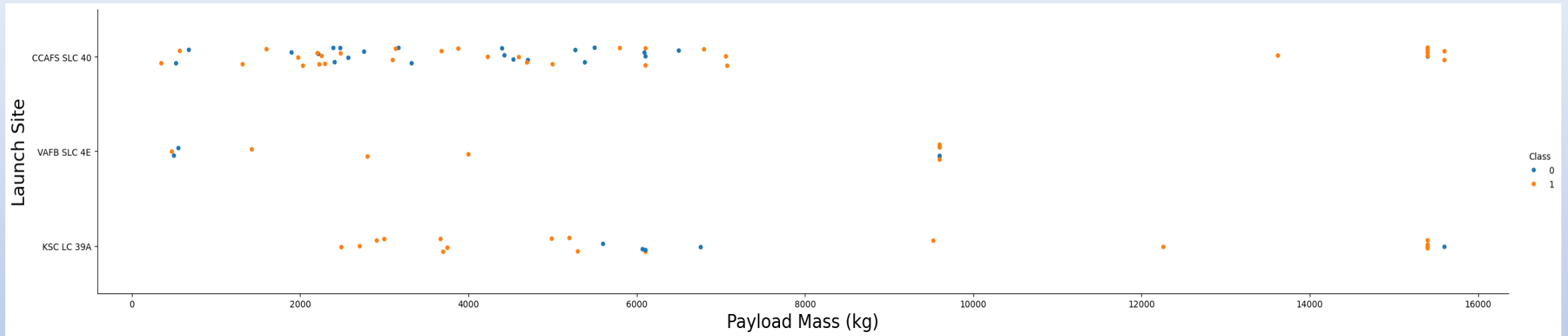


## Explanation:

- 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:

- Payload Mass Vs. Launch Site scatter point chart you will find that the VAFB-SLC launchsite there are no rockets launched for heavypayload mass(greater than 10000).

# Success Rate vs. Orbit Type

- Explanation:

Orbits with 100% success rate:

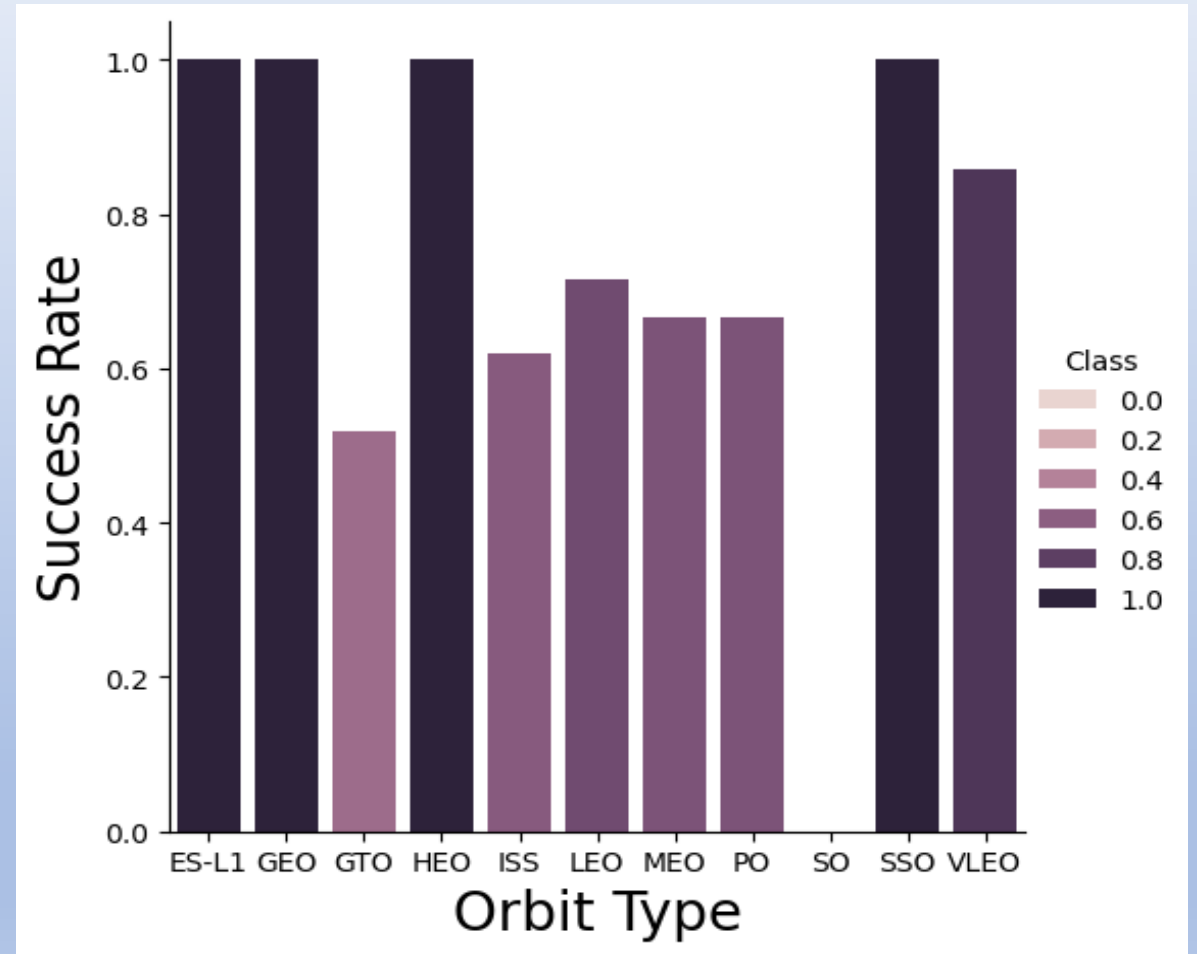
- ES-L1, GEO, HEO, SSO

Orbits with 0% success rate:

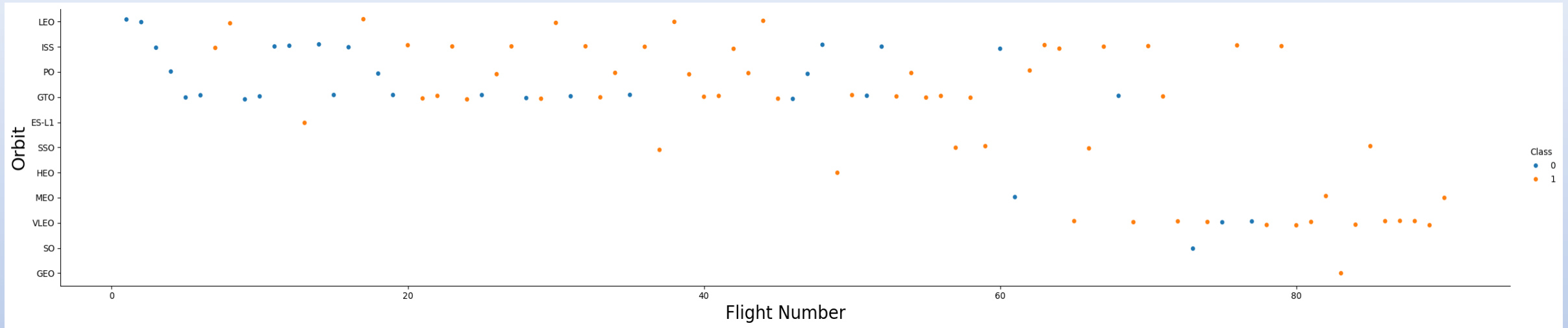
- SO

Orbits with success rate between 50% and 85%:

- GTO, ISS, LEO, MEO, PO



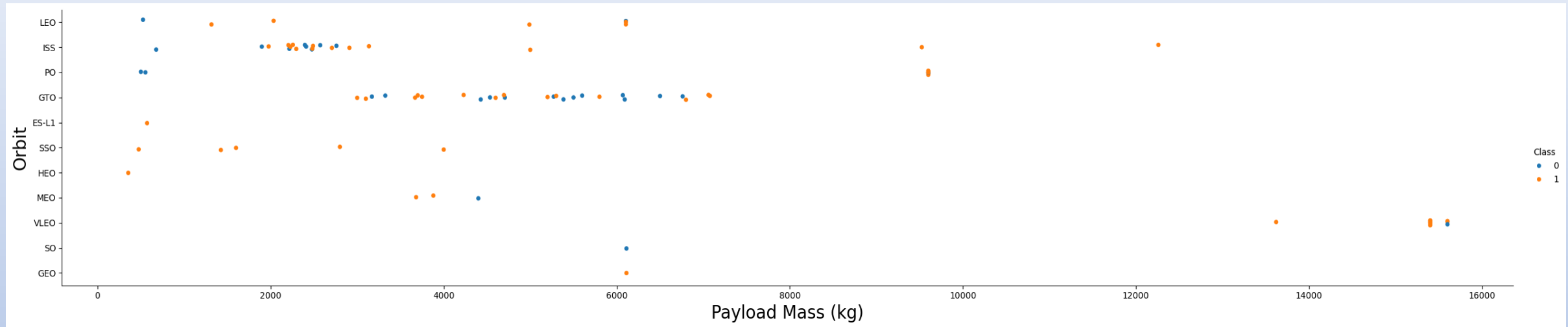
# Flight Number vs. Orbit Type



## Explanation:

- The LEO orbit, success seems to be related to the number of flights. Conversely, in the GTO orbit, there appears to be no relationship between flight number and success.

# Payload vs. Orbit Type



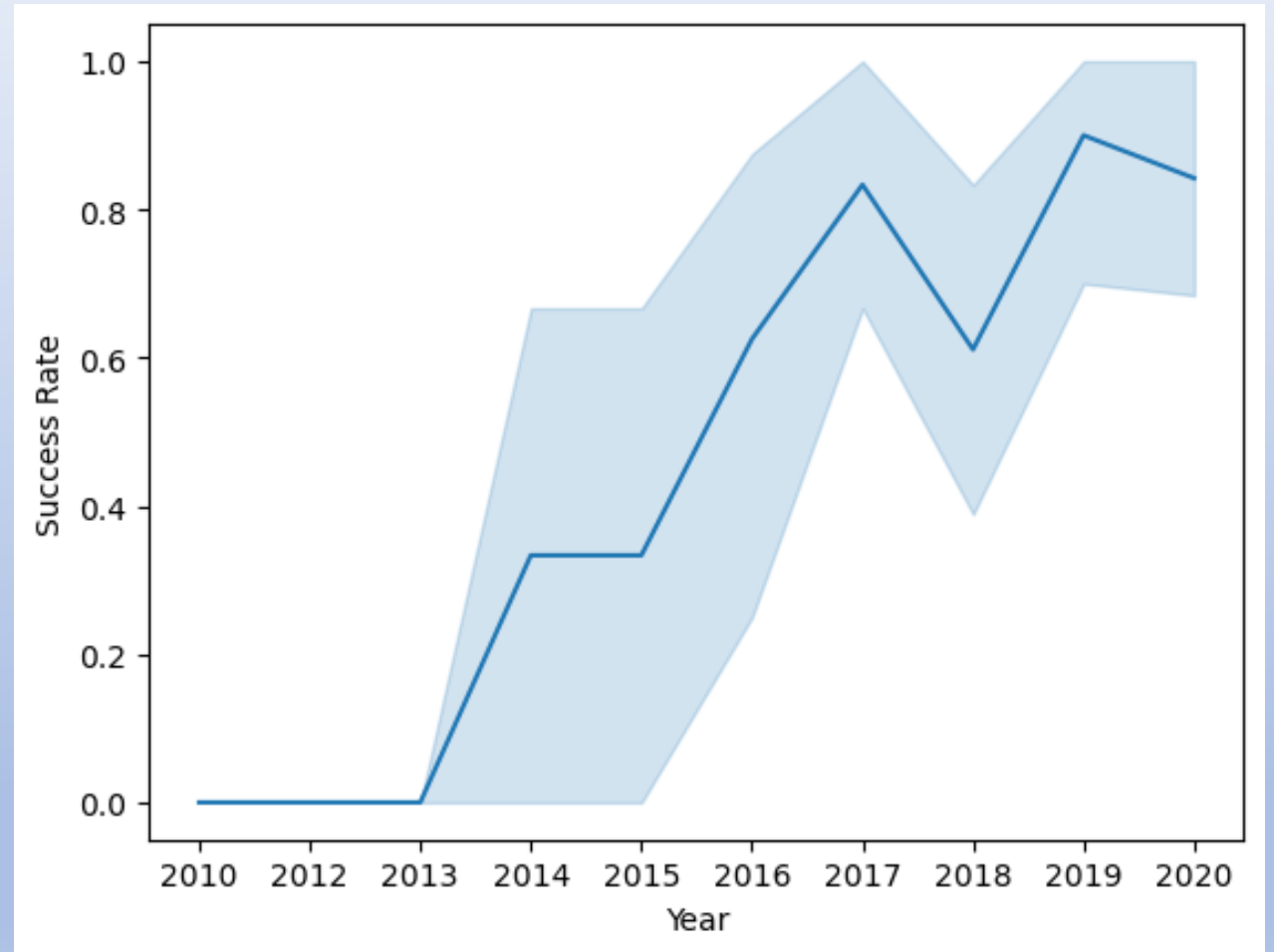
## Explanation:

- With heavy payloads the successful landing or positive landing rate are more for Polar, LEO and ISS.
- For GTO, it's difficult to distinguish between successful and unsuccessful landings as both outcomes are present.

# Launch Success Yearly Trend

## Explanation:

- The success rate since 2013 kept increasing till 2020





# All Launch Site Names

```
[10]: %sql select distinct(Launch_Site) from SPACEXTBL
```

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

```
[10]: .....
```

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

Explanation:

- Displaying the name of unique launch site using DISTINCT function

# Launch Site Names Begin with 'CCA'

```
[11]: %sql select * from SPACEXTBL where Launch_Site Like 'CCA%' limit 5;
```

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

```
[11]: .....
```

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	7:44:00	F9 v1.0 B0005	CCAFS LC-40	Dragon demo flight C2	525	LEO (ISS)	NASA (COTS)	Success	No attempt
2012-10-08	0:35:00	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 site begins with string CCA.

# Total Payload Mass

```
[12]: %sql select sum(PAYLOAD_MASS__KG_) as Total_Payload_Mass from SPACEXTBL where Customer = 'NASA (CRS)'
```

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

```
[12]: .....
```

<u>Total_Payload_Mass</u>
---------------------------

45596
-------

## Explanation:

- Calculating Total Payload Mass using SUM function

# Average Payload Mass by F9 v1.1

```
[14]: %sql select avg(payload_mass__kg_) as average_payload_mass from SPACEXTBL where booster_version like '%F9 v1.1%';
      * sqlite:///my_data1.db
      ,Done.
[14]: .....
```

<u>average_payload_mass</u>
2534.6666666666665

## Explanation:

- Displaying Average Payload Mass carried by booster version F9 v1.1
- Calculating average payload mass using AVG function

# First Successful Ground Landing Date

```
[15]: %sql select min(date) as First_Succesful_Landing from SPACEXTBL where Landing_Outcome = 'Success (ground pad)'  
      * sqlite:///my_data1.db  
      ,Done.  
[15]: .....,  
      First_Succesful_Landing  
      2015-12-22
```

## Explanation:

- Displaying the first successful landing outcome on ground pad
- SQL Min function was used to find the first successful landing date.



# Successful Drone Ship Landing with Payload between 4000 and 6000

```
[16]: %sql select Booster_Version, Payload_Mass_Kg_ from SPACEXTBL where Landing_Outcome = 'Success (ground pad)' and Payload_Mass_Kg_ between 4000 and 6000
* sqlite:///my_data1.db
,Done.
```

```
[16]: .....
```

Booster_Version	PAYLOAD_MASS_KG_
F9 FT B1032.1	5300
F9 B4 B1040.1	4990
F9 B4 B1043.1	5000

## Explanation:

- Displaying the names of boosters which have successfully landed on drone ship and had payload mass greater than 4000 but less than 6000
- SQL Where Clause was used to find the names of boosters.

# Total Number of Successful and Failure Mission Outcomes

```
[17]: %sql select Mission_Outcome, count(*) as Total_Number from SPACEXTBL group by Mission_Outcome;  
* sqlite:///my_data1.db  
,Done.
```

```
[17]: .....
```

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

## Explanation:

- Displaying the count of Failed and Success Mission
- SQL Count function is used in finding the result

# Boosters Carried Maximum Payload

```
[18]: %sql select Booster_Version from SPACEXTBL where Payload_Mass__KG_ = (select max(Payload_Mass__KG_) from SPACEXTBL)
* sqlite:///my_data1.db
,Done.
```

```
[18]: .....
```

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 Booster Version that carried maximum payload mass using SQL subquery.

# 2015 Launch Records

```
•[30]: %sql select substr(Date, 6,2) as Monthname, Date, Booster_Version, Launch_Site, Landing_Outcome from SPACEXTBL where Landing_Outcome = 'Failure (drone ship)'
* sqlite:///my_data1.db
,Done.
```

```
[30]: .....
```

Monthname	Date	Booster_Version	Launch_Site	Landing_Outcome
01	2015-01-10	F9 v1.1 B1012	CCAFS LC-40	Failure (drone ship)
04	2015-04-14	F9 v1.1 B1015	CCAFS LC-40	Failure (drone ship)

## Explanation:

- 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.

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

```
[38]: %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
```

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

```
[38]: .....
```

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

## Explanation:

- Ranking 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.

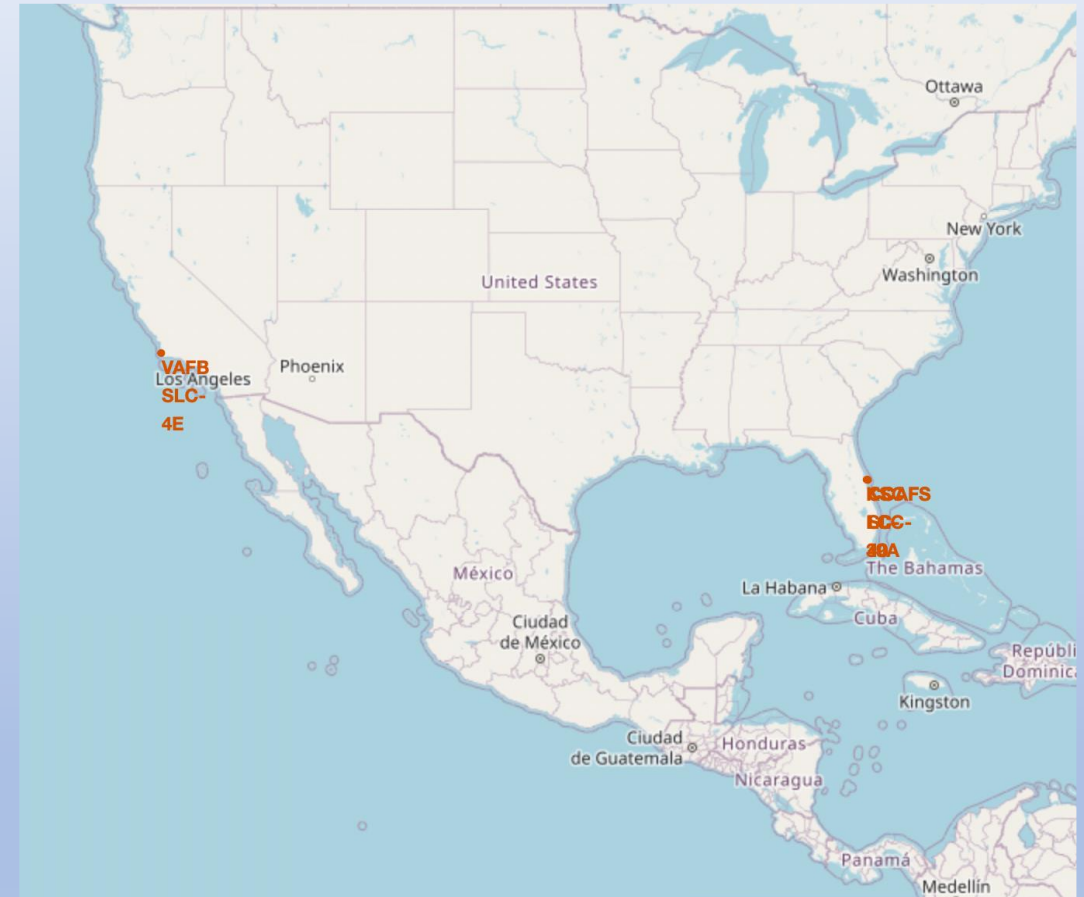
Section 3

# Launch Sites Proximities Analysis

# Launch Sites Location Markers on 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 minimizes the risk of having any debris dropping or exploding near people.

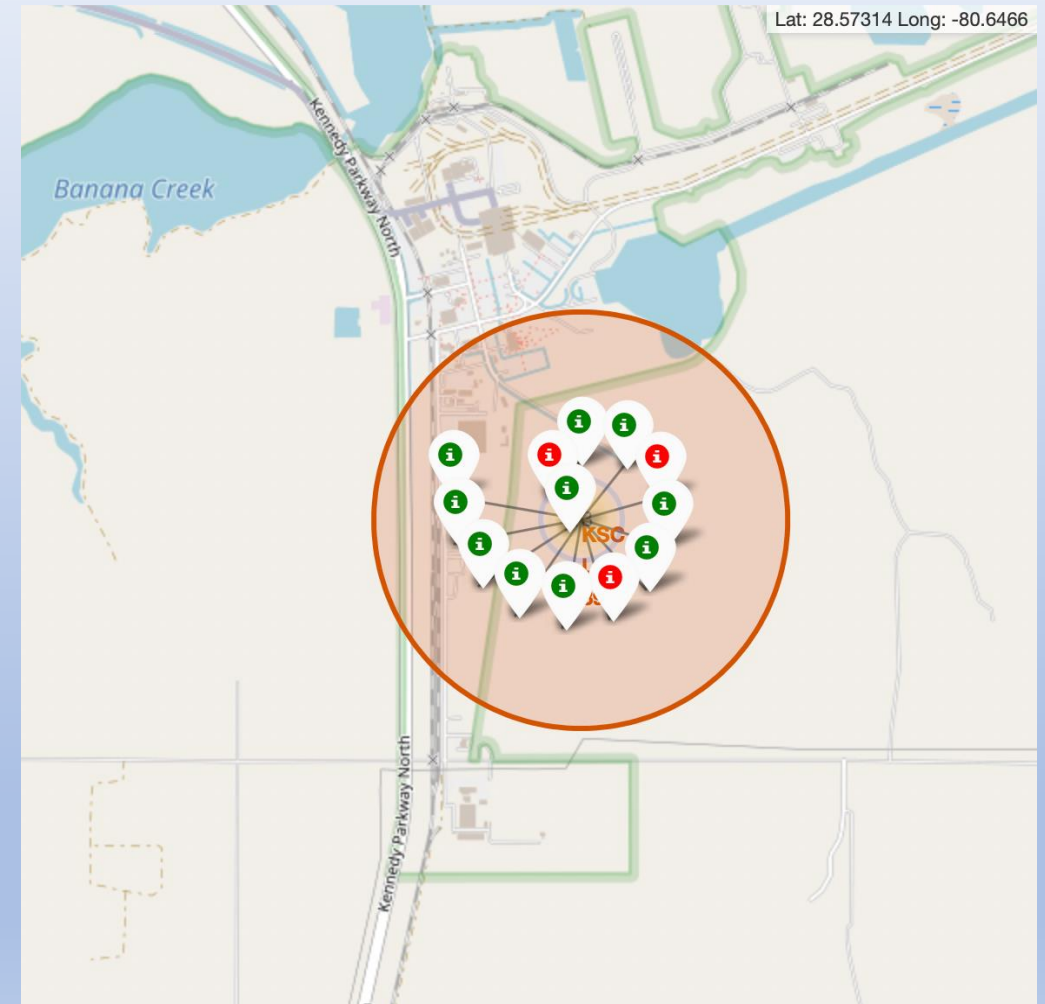




# Colored Labeled Launch Records on the Map

## Explanation:

- From the color-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.

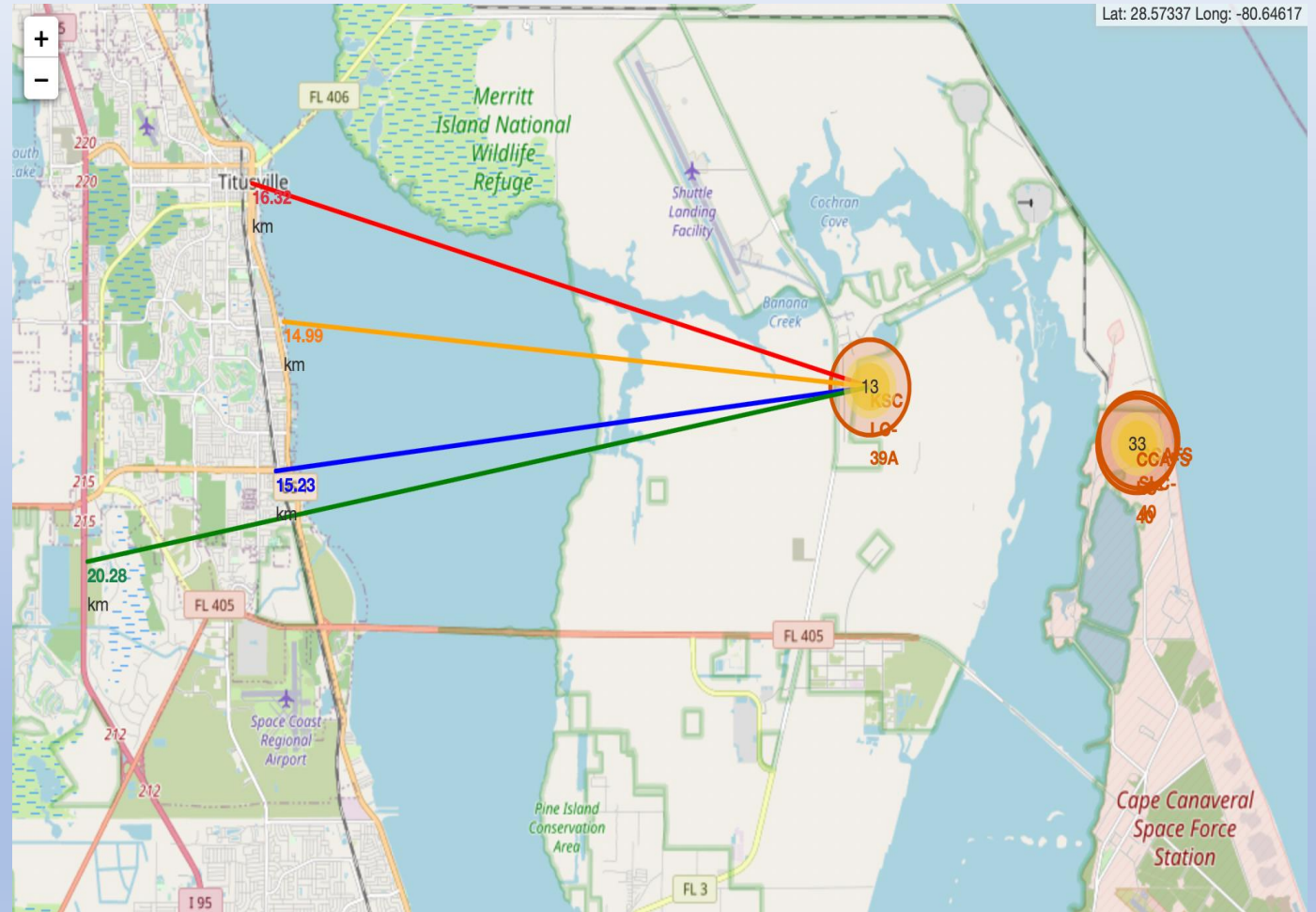




# Launch Site KSC LC-39A proximity distance

## Explanation:

- From the visual analysis of the launch site KSC LC-39A we can clearly see that it is:
  - relatively close to railway (15.23 km)
  - relatively close to highway (20.28 km)
  - relatively close to coastline (14.99 km)
- The launch site KSC LC-39A is relatively 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.

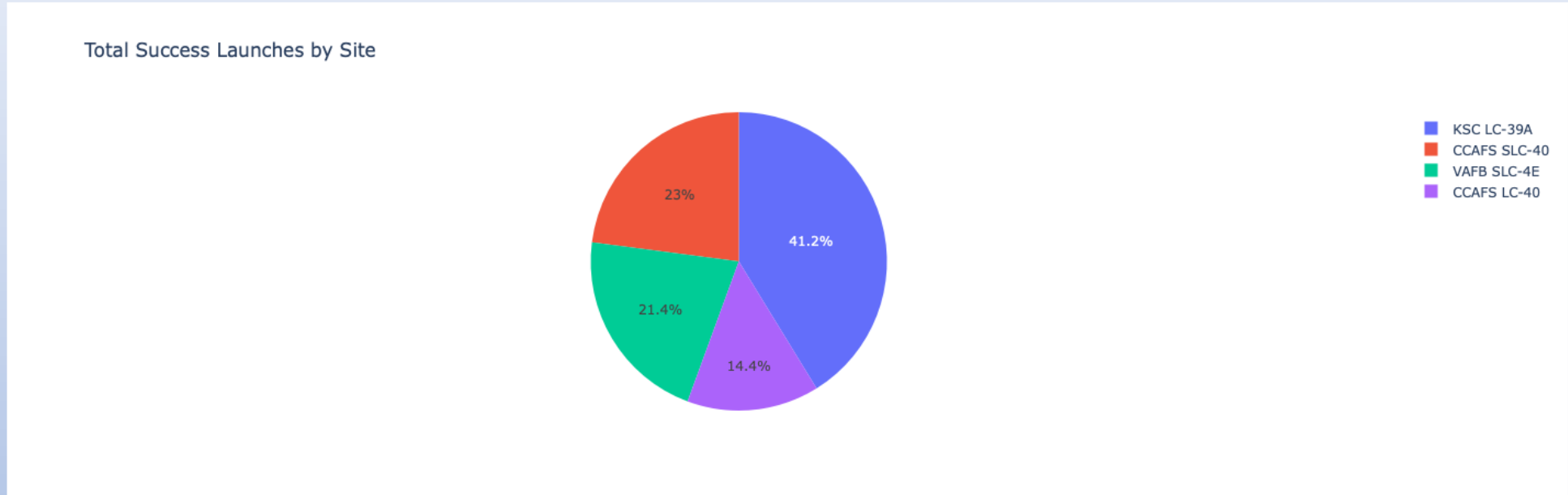


The background of the slide is a close-up, artistic photograph of a printed circuit board (PCB). The board is dark, and the intricate circuit traces are highlighted in a vibrant, glowing red. Numerous small, cylindrical electronic components, likely capacitors or resistors, are visible, some of which also appear to be glowing with a warm, orange-red light. The overall aesthetic is high-tech and digital, suggesting themes of data, technology, and engineering.

Section 4

# Build a Dashboard with Plotly Dash

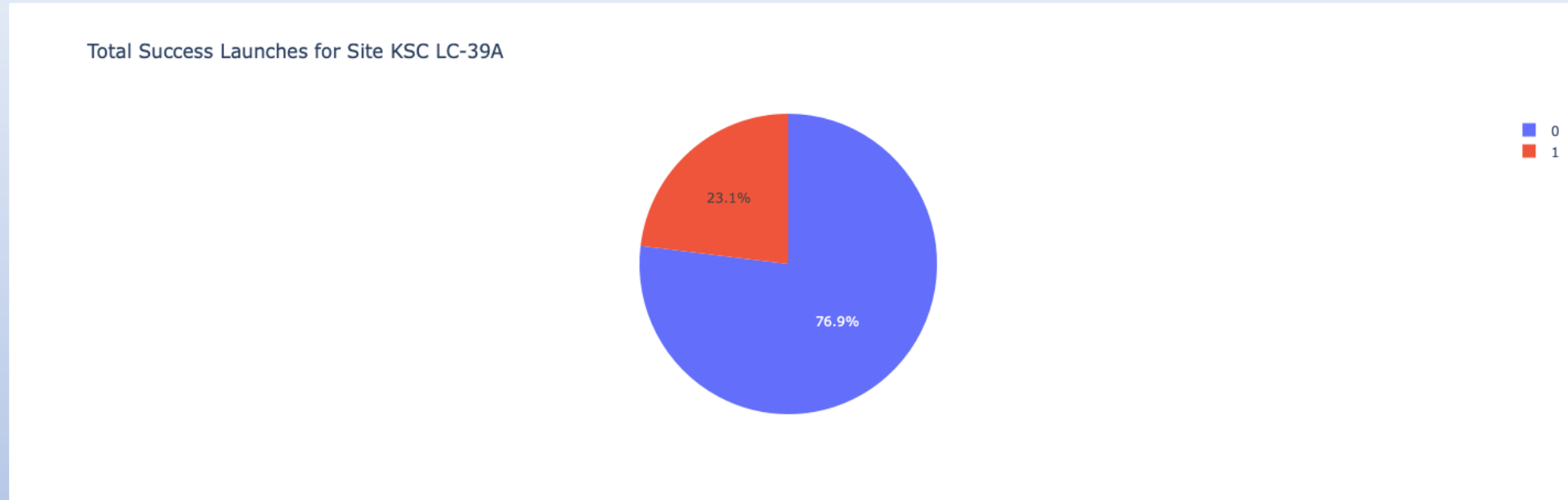
# Success Rate for all Launch Sites



## 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.



# Pay Load Mass vs. Launch Outcome for All Sites

## Explanation:

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



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 test set

	LogReg	SVM	Tree	KNN
<b>Jaccard_Score</b>	0.800000	0.800000	0.800000	0.800000
<b>F1_Score</b>	0.888889	0.888889	0.888889	0.888889
<b>Accuracy</b>	0.833333	0.833333	0.833333	0.833333

## Scores and Accuracy of the data set

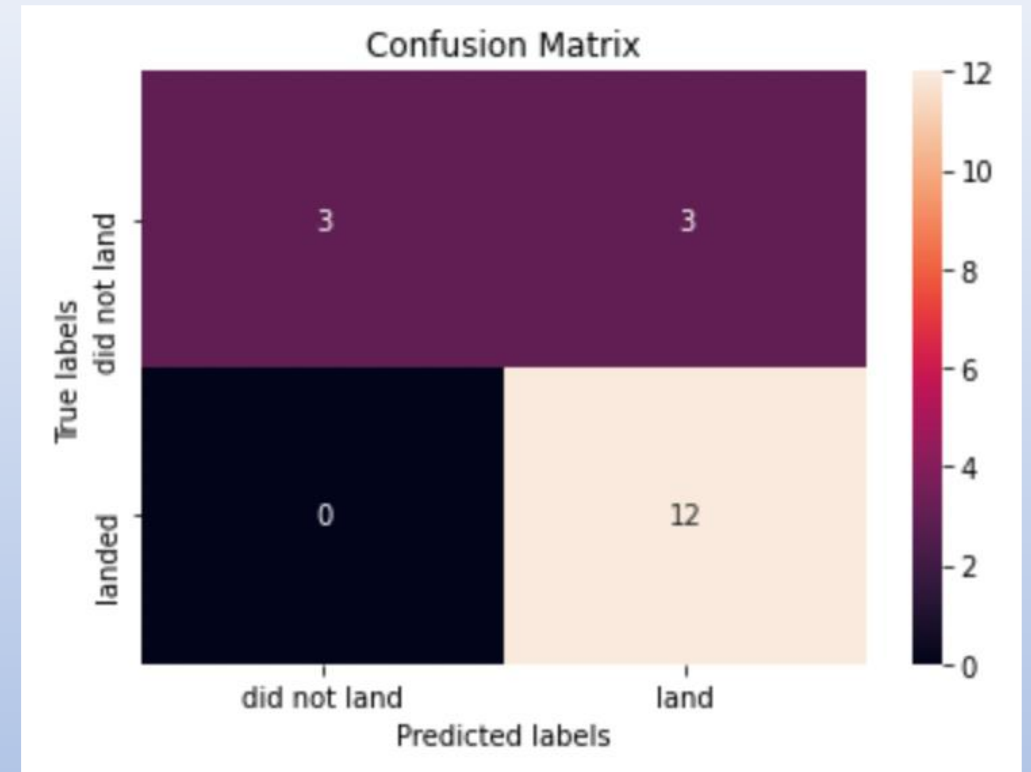
	LogReg	SVM	Tree	KNN
<b>Jaccard_Score</b>	0.833333	0.845070	0.882353	0.819444
<b>F1_Score</b>	0.909091	0.916031	0.937500	0.900763
<b>Accuracy</b>	0.866667	0.877778	0.911111	0.855556

# 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. • Orbits ES-L1, GEO, HEO and SSO have 100% success rate.

# Appendix

Special Thanks to:

[Instructor](#)

[Coursera](#)

[IBM](#)

Thank you!

