

# Data Science Capstone Project

Maryam Beik

29/03/2024



# Outline

---



• Executive Summary.....	3
• Introduction.....	5
• Methodology.....	6
• Results.....	17
• Conclusion.....	47
• Appendix.....	48

# Executive Summary

---

- Methodologies

## Data -

Data source - public SpaceX API data and Wikipedia

Exploratory Data Analysis - Python libraries, SQL, visualization tools, folium maps, and dashboards

Feature Engineering - Extracting relevant columns, Converting categorical variables into binary format using one-hot encoding and utilizing GridSearchCV to identify the optimal parameters for machine learning models.

Data cleaning and preprocessing - Data standardization , dealing with missing values and outliers.

# Executive Summary

---

- Methodologies

## Machine Learning Algorithms -

- Logistic Regression
- Support Vector Machine
- Decision Tree Classifier
- K Nearest Neighbour

Based on the training of the mentioned algorithms, it can be reasonably conclude that the results achieved an accuracy rate of approximately 83.33%.

That is for sure more data is needed for better model determination and accuracy.

# Introduction

---

- **Project background and context**

- The Space Exploration Technologies (SpaceX) mission is to make affordable rockets. Its first two rockets were the Falcon 1 (first launched in 2006) and the larger Falcon 9 (first launched in 2010), which were designed to cost much less than competing rockets. ( \$62 million vs. \$165 million USD)
- Space X has largely gained an advantage due to its capability to recover parts of the rocket, particularly Stage 1.
- **Y Space Y has been competing with Space X.**

- **Problems you want to find answers**

- **Space Y wants us to develop a predictive model to predict successful Stage 1 recovery.**

Section 1

# Methodology

# Methodology

---

- Data collection methodology:
  - Data collection is conducted from SpaceX API and Web scraping from Space X's Wikipedia
- Perform data wrangling
  - Classifying true landing as **successful** and **unsuccessful**
- Perform exploratory data analysis (EDA) using visualization and SQL
- Perform interactive visual analytics using Folium and Plotly Dash
- Perform predictive analysis using classification models
  - Tuned for best parameters using GridSearchCV

# Data Collection

---

## Overview

Data collection process involved a combination of API requests from Space X public API and web scraping data from a table in Space X's Wikipedia entry.

### Space X API Data Columns:

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

### Wikipedia Webscrape Data Columns:

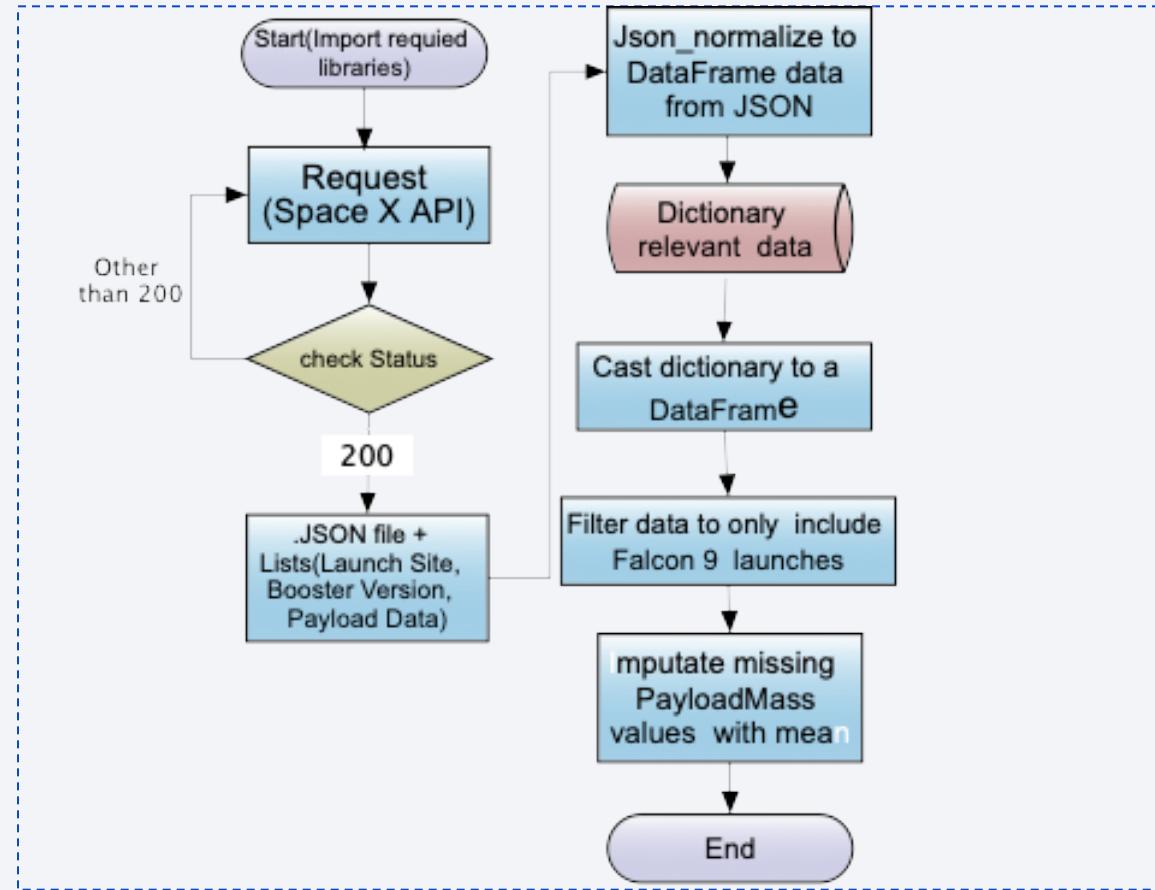
Flight No., Launch site, Payload, PayloadMass, Orbit, Customer, Launch outcome, Version, Booster, Booster landing, Date, Time

# Data Collection – SpaceX API

- Data Collection - SpaceX API

- GitHub URL:

[https://github.com/MaryamBeik/testrepo/blob/main/SpaceX\\_Machine%20Learning%20Prediction\\_Part\\_5.ipynb](https://github.com/MaryamBeik/testrepo/blob/main/SpaceX_Machine%20Learning%20Prediction_Part_5.ipynb)

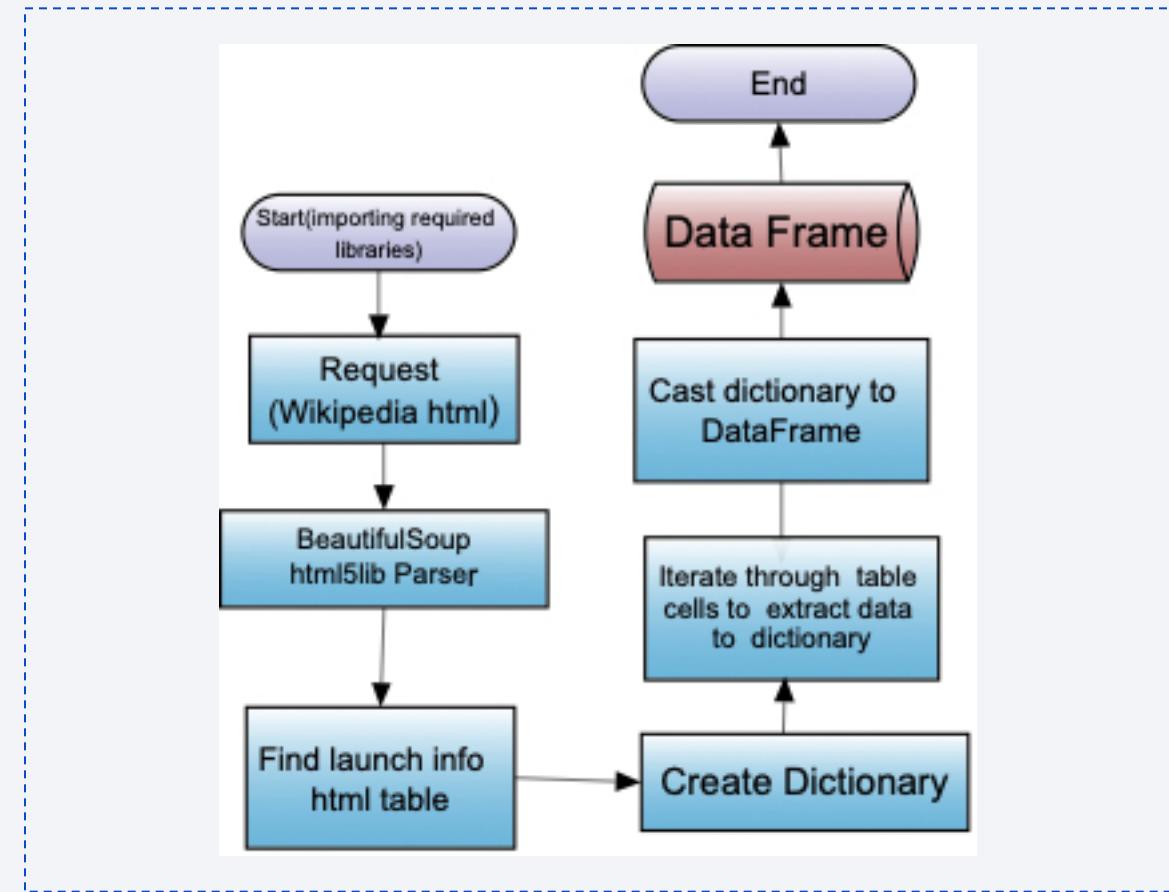


# Data Collection – Web Scraping

- Data Collection - Web Scraping

- GitHub URL:

<https://github.com/MaryamBeik/testrepo/blob/main/jupyter-labs-webscraping.ipynb>



# Data Wrangling

---

Creating a training label (Y) with landing outcomes where successful = 1 & failure = 0.

Outcome column has two components: ‘Mission Outcome’ ‘Landing Location’

New training label column ‘class’ with a value of 1 if ‘Mission Outcome’ is True and 0 otherwise.

## Value Mapping:

True ASDS, True RTLS, & True Ocean – set to -> 1

None None, False ASDS, None ASDS, False Ocean, False RTLS – set to ->0

GitHub URL:<https://github.com/MaryamBeik/testrepo/blob/main/labs-jupyter-spacex-Data%20wrangling.ipynb>

# EDA with Data Visualization

---

Exploratory Data Analysis performed on

- **Variables:**

Flight Number, Payload Mass, Launch Site, Orbit, Class and Year.

- **Plots :**

Scatter plots, line charts, and bar plots

to compare relationships between variables to decide if a relationship exists so that they could be used in training the machine learning model

Flight Number vs. Payload Mass, Flight Number vs. Launch Site, Payload Mass vs. Launch Site, Orbit vs. Success Rate, Flight Number vs. Orbit, Payload vs Orbit, and Success Yearly Trend

GitHub URL:<https://github.com/MaryamBeik/testrepo/blob/main/jupyter-labs-eda-dataviz.ipynb>

## EDA with SQL

---

- Loaded data set into IBM DB2 Database
- Queried using SQL Python integration
- Queries were made to get a better understanding of the dataset.
- Queried information about launch site names, mission outcomes, various pay load sizes of customers and booster versions, and landing outcomes
- GitHub URL: [https://github.com/MaryamBeik/testrepo/blob/main/jupyter-labs-eda-sql-coursera\\_sqlite.ipynb](https://github.com/MaryamBeik/testrepo/blob/main/jupyter-labs-eda-sql-coursera_sqlite.ipynb)

# Build an Interactive Map with Folium

---

- **Map objects**

## **markers, circles, lines, circleMarker, markerCluster**

Folium maps mark Launch Sites, successful and unsuccessful landings, and a proximity example to key locations: Railway, Highway, Coast, and City.

This allows us to understand why launch sites may be located where they are. Also visualizes successful landings relative to location.

- GitHub URL : <https://github.com/MaryamBeik/testrepo/blob/main/Interactive%20Visual%20Analytics%20with%20Folium.ipynb>

# Build a Dashboard with Plotly Dash

---

## Plots in Dashboard

pie chart, scatter plot

Pie chart- to show distribution of successful landings across all launch sites and also show individual launch site success rates.

Scatter plot with two inputs- All sites or individual site and payload mass on a slider between 0 and 10000 kg.

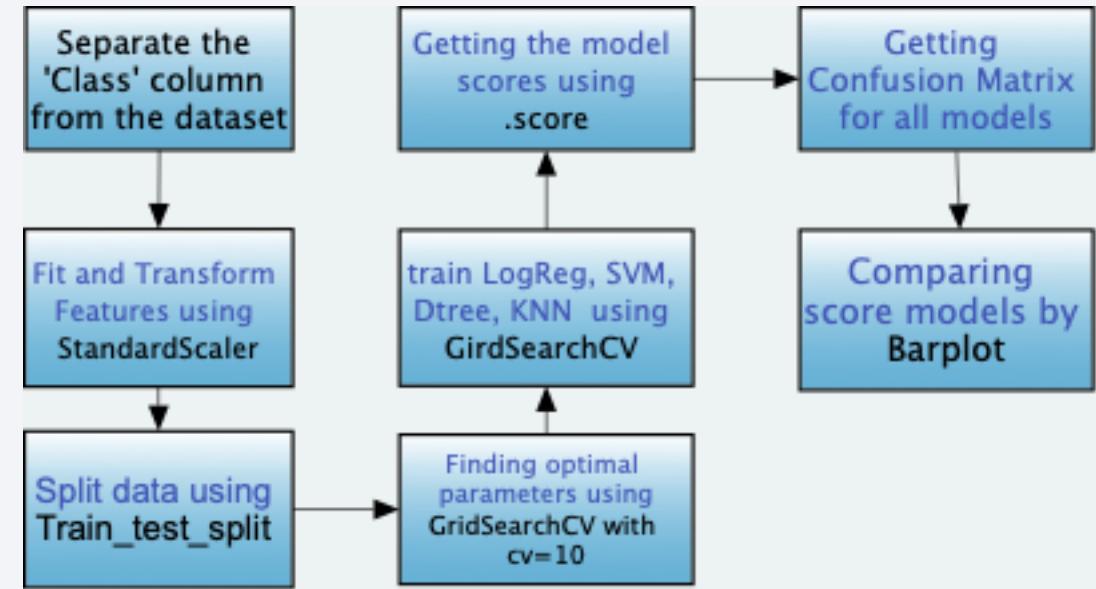
The pie chart is used to visualize launch site success rate.

The scatter plot can help us see how success varies across launch sites, payload mass, and booster version category.

- GitHub URL: <https://github.com/MaryamBeik/testrepo/blob/main/Create%20Dashboard%20using%20plotly%20and%20Dash.ipynb>

# Predictive Analysis (Classification)

---

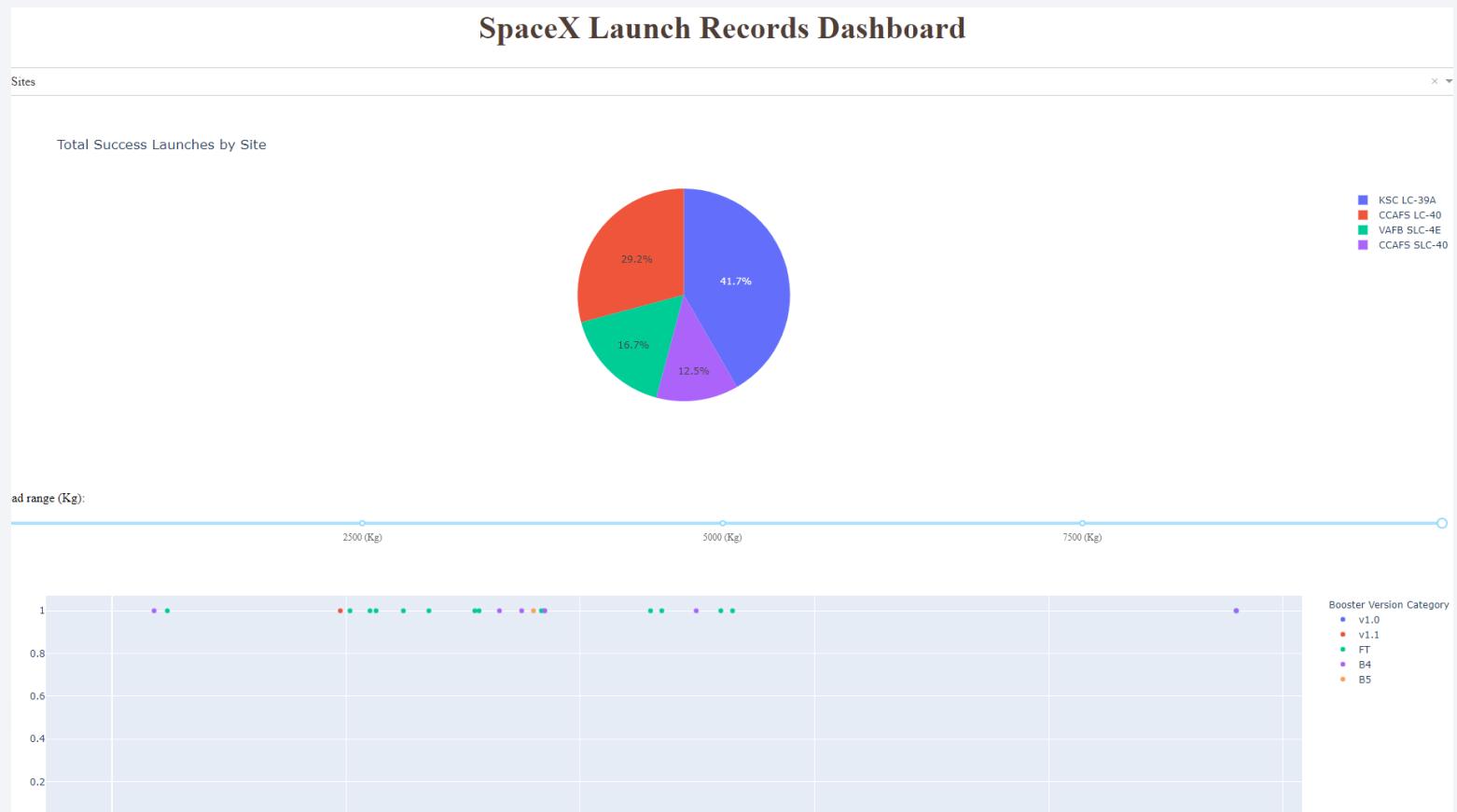


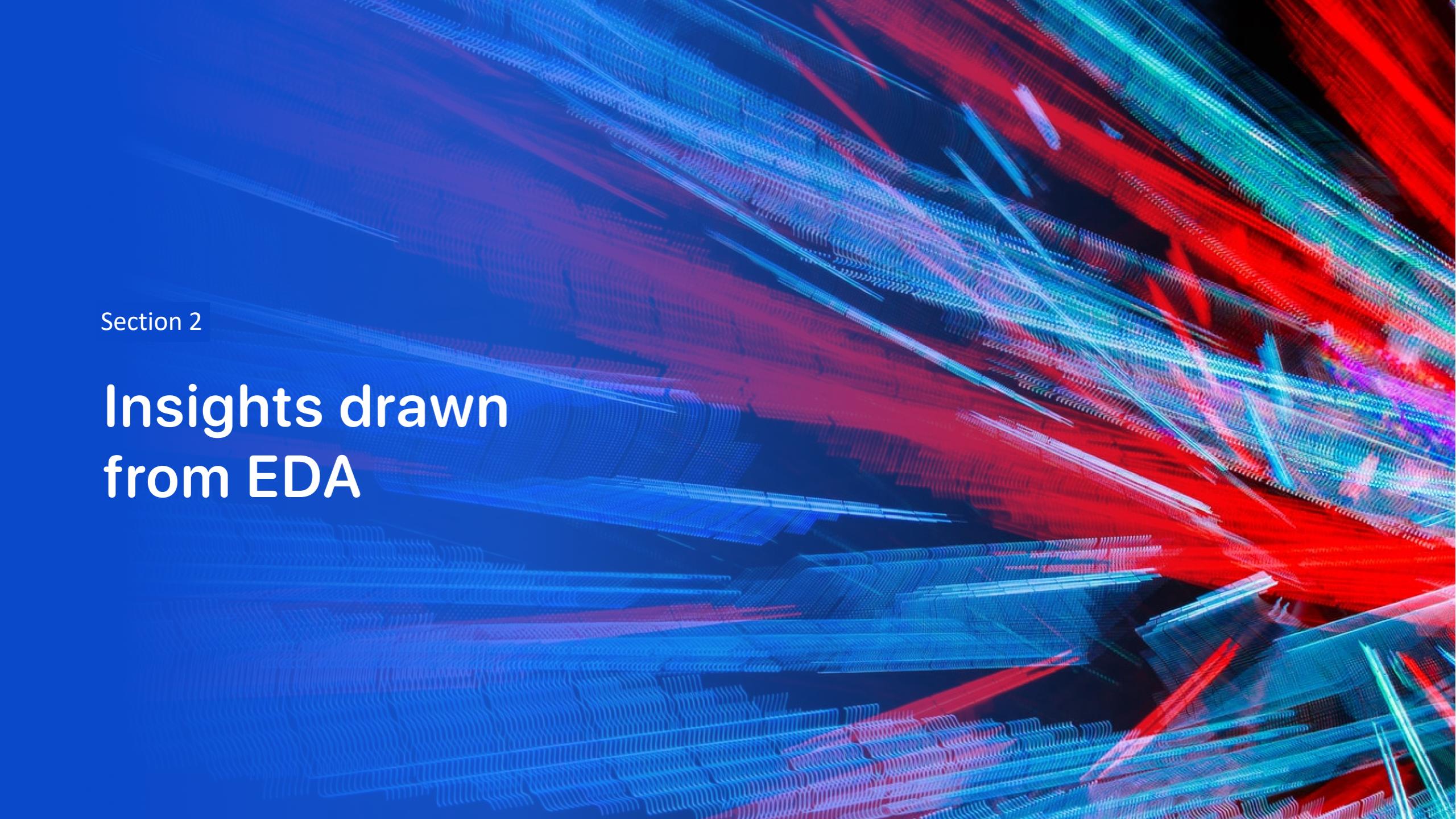
- GitHub URL: [https://github.com/MaryamBeik/testrepo/blob/main/SpaceX\\_Machine%20Learning%20Prediction\\_Part\\_5.ipynb](https://github.com/MaryamBeik/testrepo/blob/main/SpaceX_Machine%20Learning%20Prediction_Part_5.ipynb)

# Results

---

This is a preview of the Plotly dashboard. The following slides will show the results of EDA with visualization, EDA with SQL, Interactive Map with Folium, and finally the results of our model with about 83% accuracy.

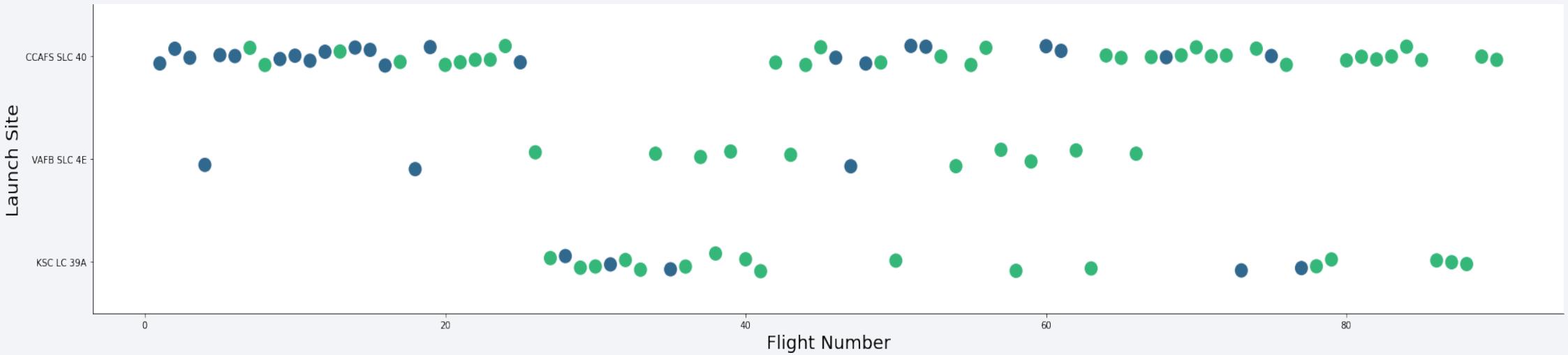


The background of the slide features a complex, abstract digital visualization. It consists of numerous thin, glowing lines that create a sense of depth and motion. The lines are primarily blue and red, with some green and white highlights. They form a grid-like structure that is more dense and vibrant towards the right side of the frame, while appearing more sparse and blue-tinted on the left. The overall effect is reminiscent of a high-energy particle simulation or a futuristic circuit board.

Section 2

## Insights drawn from EDA

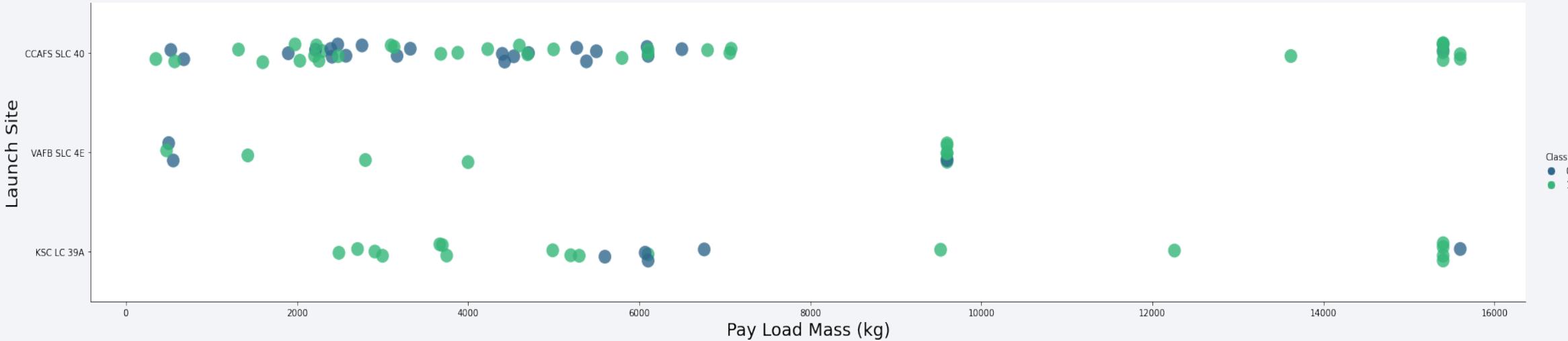
# Flight Number vs. Launch Site



Green indicates successful launch; Purple indicates unsuccessful launch.

Graphic suggests an increase in success rate over time (indicated in Flight Number). Likely a big breakthrough around flight 20 which significantly increased success rate. CCAFS appears to be the main launch site as it has the most volume.

# Payload vs. Launch Site

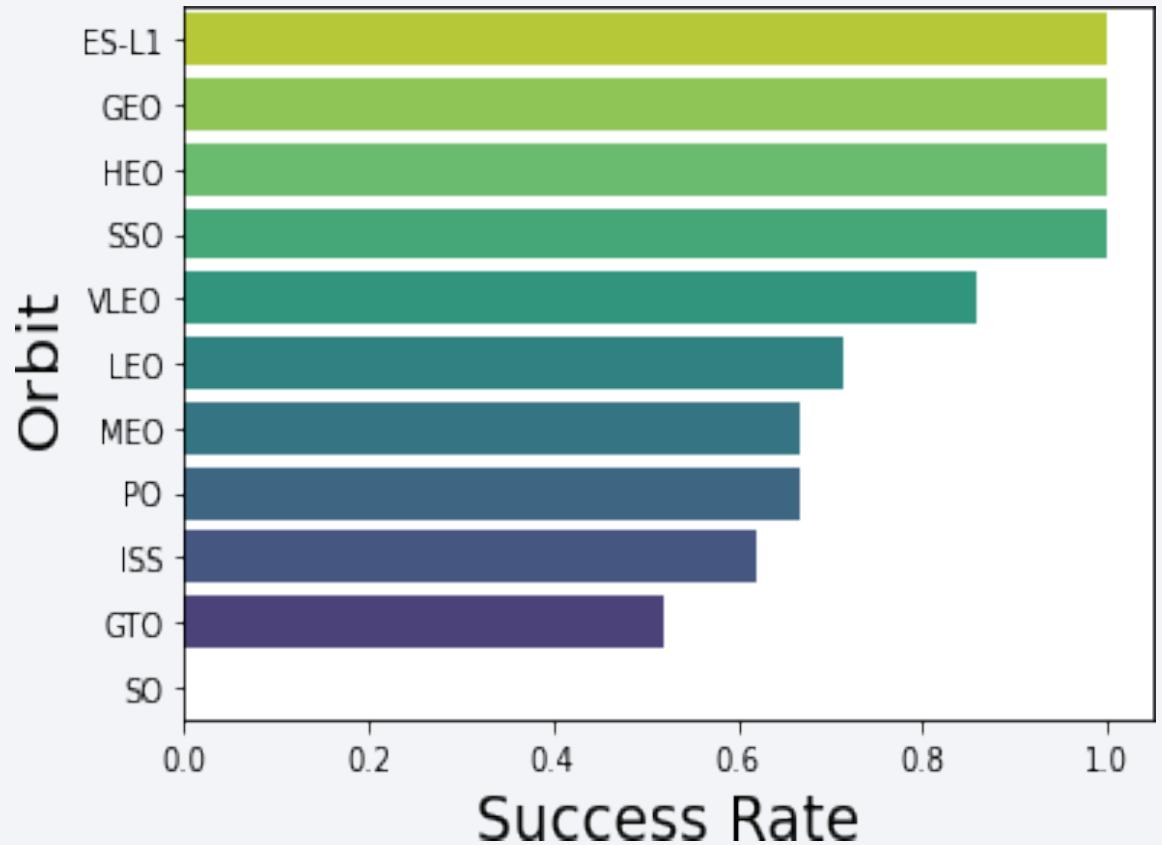


Green indicates successful launch; Purple indicates unsuccessful launch.

Payload mass appears to fall mostly between 0-6000 kg. Different launch sites also seem to use different payload mass.

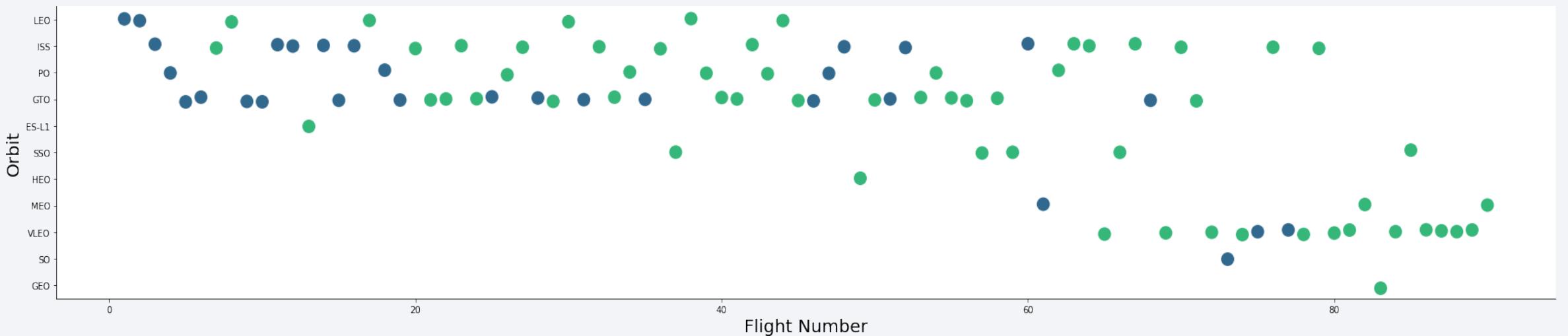
# Success Rate vs. Orbit Type

ES-L1 (1), GEO (1), HEO (1) have 100% success rate (sample sizes in parenthesis) SSO (5) has 100% success rate  
VLEO (14) has decent success rate and attempts  
SO (1) has 0% success rate  
GTO (27) has the around 50% success rate but largest sample



Success Rate Scale with 0 as 0%  
0.6 as 60% 1 as 100%

# Flight Number vs. Orbit Type

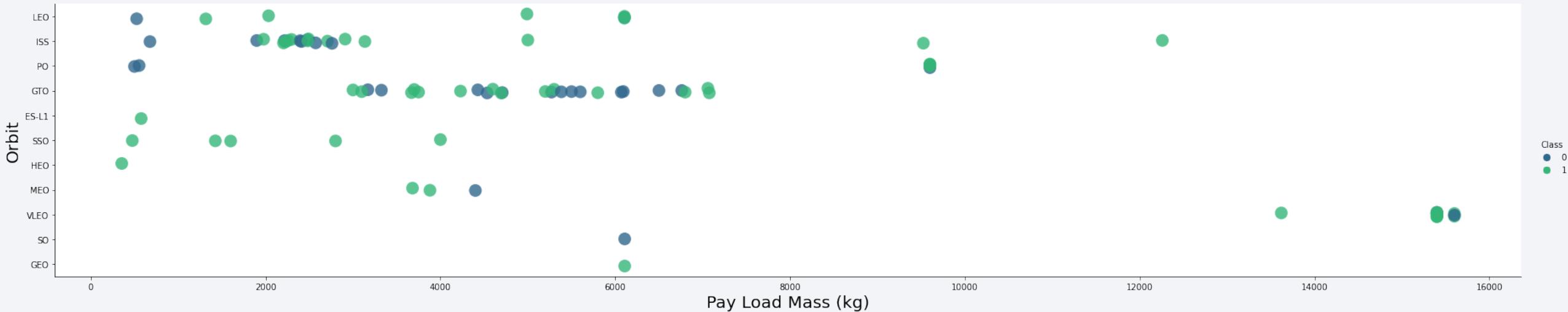


Green indicates successful launch; Purple indicates unsuccessful launch.

Launch Orbit preferences changed over Flight Number. Launch Outcome seems to correlate with this preference.

SpaceX started with LEO orbits which saw moderate success LEO and returned to VLEO in recent launches. SpaceX appears to perform better in lower orbits or Sun-synchronous orbits

# Payload vs. Orbit Type



Green indicates successful launch; Purple indicates unsuccessful launch.

Payload mass seems to correlate with orbit

LEO and SSO seem to have relatively low payload mass

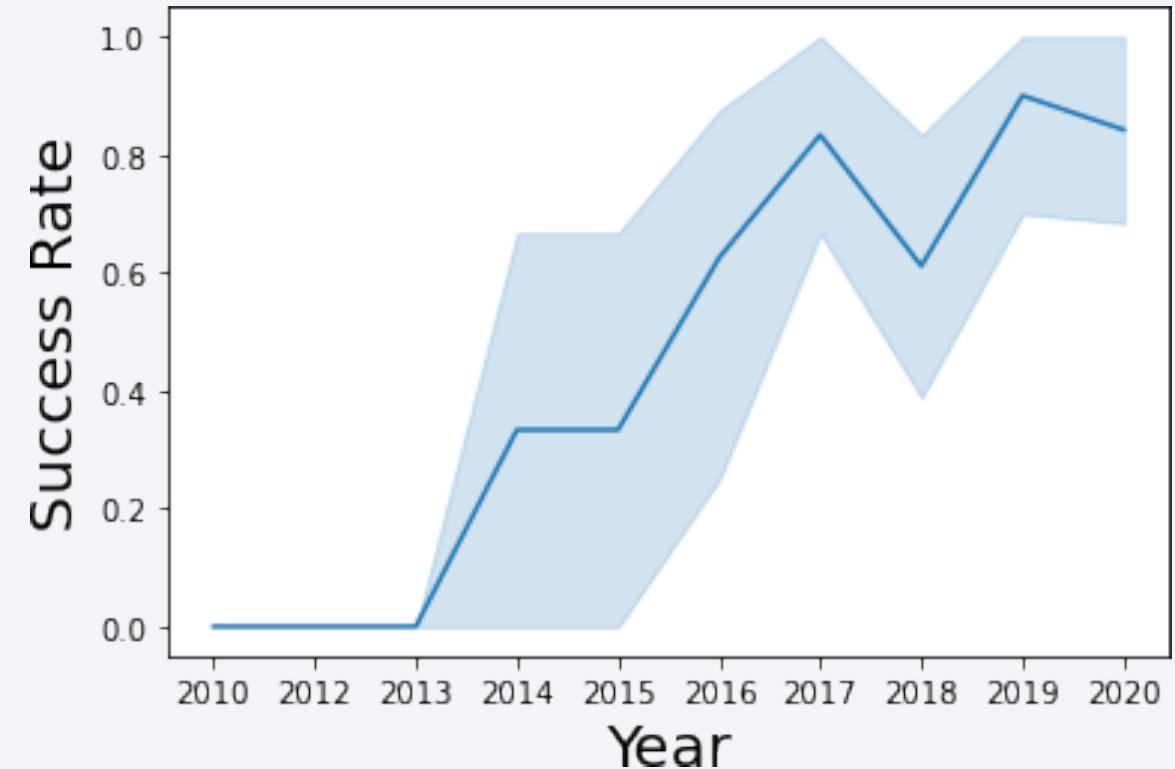
The other most successful orbit VLEO only has payload mass values in the higher end of the range

# Launch Success Yearly Trend

---

Success generally increases over time since 2013 with a slight dip in 2018

Success in recent years at around 80%



95% confidence interval  
(light blue shading)

# EDA with SQL

---

**Exploratory Data Analysis**

**SQL DB2**

# All Launch Site Names

---

Query unique launch site names from database.

In [4]: `%%sql  
SELECT UNIQUE LAUNCH_SITE  
FROM SPACEXDATASET;  
* ibm_db_sa://ftb12020:***@0c77d6f:  
Done.`

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

CCAFS SLC-40 and CCAFSSLC-40 likely all represent the same launch site with data entry errors. CCAFS LC-40 was the previous name.  
Likely only 3 unique launch\_site values:  
CCAFS SLC-40, KSC  
LC-39A, VAFB SLC-4E

# Launch Site Names Begin with 'CCA'

First five entries in database with Launch Site name beginning with CCA.

In [5]:

```
%%sql
SELECT *
FROM SPACEXDATASET
WHERE LAUNCH_SITE LIKE 'CCA%'
LIMIT 5;
```

\* ibm\_db\_sa://ftb12020:\*\*\*@0c77d6f2-5da9-48a9-81f8-86b520b87518.bs2io90l08kqb1od8lcg.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

---

## Total payload carried by boosters from NASA

```
%%sql
SELECT SUM(PAYLOAD_MASS__KG_) AS SUM_PAYLOAD_MASS_KG
FROM SPACEXDATASET
WHERE CUSTOMER = 'NASA (CRS)';

* ibm_db_sa://ftb12020:***@0c77d6f2-5da9-48a9-81f8-86
Done.
```

sum_payload_mass_kg
45596

This query sums the total payload mass in kg where NASA was the customer.

CRS stands for Commercial Resupply Services which indicates that these payloads were sent to the International Space Station (ISS).

# Average Payload Mass by F9 v1.1

## Average payload mass carried by booster version F9 v1.1

```
%%sql
SELECT AVG(PAYLOAD_MASS__KG_) AS AVG_PAYLOAD_MASS_KG
FROM SPACEXDATASET
WHERE booster_version = 'F9 v1.1'
```

```
* ibm_db_sa://ftb12020:***@0c77d6f2-5da9-48a9-81f8-8e
```

```
Done.
```

avg_payload_mass_kg
2928

This query calculates the average payload mass of launches which used booster version F9 v1.1

Average payload mass of F9 1.1 is on the low end of our payload mass range

# First Successful Ground Landing Date

---

## Dates of the first successful landing outcome on ground pad

```
%%sql
SELECT MIN(DATE) AS FIRST_SUCCESS
FROM SPACEXDATASET
WHERE landing_outcome = 'Success (ground pad)';
* ibm_db_sa://ftb12020:***@0c77d6f2-5da9-48a9-81
Done.
```

first_success
2015-12-22

This query returns the first successful ground pad landing date.

First ground pad landing wasn't until the end of 2015.

Successful landings in general appear starting 2014.

# Successful Drone Ship Landing with Payload between 4000 and 6000

**Names of boosters which have successfully landed on drone ship and had payload mass greater than 4000 but less than 6000**

```
%%sql
SELECT booster_version
FROM SPACEXDATASET
WHERE landing_outcome = 'Success (drone ship)' AND payload_mass_kg_ BETWEEN 4001 AND 5999;
* ibm_db_sa://ftb12020:***@0c77d6f2-5da9-48a9-81f8-86b520b87518.bs2io90l08kqb1od81cg.database.
Done.
```

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

This query returns the four booster versions that had successful drone ship landings and a payload mass between 4000 and 6000 noninclusively.

# Total Number of Successful and Failure Mission Outcomes

---

## Total number of successful and failure mission outcomes

```
%%sql
SELECT mission_outcome, COUNT(*) AS no_outcome
FROM SPACEXDATASET
GROUP BY mission_outcome;
* ibm_db_sa://ftb12020:***@0c77d6f2-5da9-48a9-8
Done.
```

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

This query returns a count of each mission outcome.

SpaceX appears to achieve its mission outcome nearly 99% of the time.

This means that most of the landing failures are intended.

Interestingly, one launch has an unclear payload status and unfortunately one failed in flight.

# Boosters Carried Maximum Payload

## Names of the booster which have carried the maximum payload mass

```
%%sql
SELECT booster_version, PAYLOAD_MASS__KG_
FROM SPACEXDATASET
WHERE PAYLOAD_MASS__KG_ = (SELECT MAX(PAYLOAD_MASS__KG_) FROM SPACEXDATASET);

* ibm_db_sa://ftb12020:***@0c77d6f2-5da9-48a9-81f8-86b520b87518.bs2io90l08kqb1
Done.
```

booster_version	payload_mass__kg_
F9 B5 B1048.4	15600
F9 B5 B1049.4	15600
F9 B5 B1051.3	15600
F9 B5 B1056.4	15600
F9 B5 B1048.5	15600
F9 B5 B1051.4	15600
F9 B5 B1049.5	15600
F9 B5 B1060.2	15600
F9 B5 B1058.3	15600
F9 B5 B1051.6	15600
F9 B5 B1060.3	15600
F9 B5 B1049.7	15600

This query returns the booster versions that carried the highest payload mass of 15600 kg.

These booster versions are very similar and all are of the F9 B5 B10xx.x variety.

This likely indicates payload mass correlates with the booster version that is used.

# 2015 Launch Records

## Failed landing\_outcomes in drone ship, booster versions, launch site names in year 2015

```
%%sql
SELECT MONTHNAME(DATE) AS MONTH, landing_outcome, booster_version, PAYLOAD_MASS_KG_, launch_site
FROM SPACEXDATASET
WHERE landing_outcome = 'Failure (drone ship)' AND YEAR(DATE) = 2015;
```

```
* ibm_db_sa://ftb12020:***@0c77d6f2-5da9-48a9-81f8-86b520b87518.bs2io90l08kqb1od8lcg.databases.app
Done.
```

MONTH	landing_outcome	booster_version	payload_mass_kg	launch_site
January	Failure (drone ship)	F9 v1.1 B1012	2395	CCAFS LC-40
April	Failure (drone ship)	F9 v1.1 B1015	1898	CCAFS LC-40

This query returns the Month, Landing Outcome, Booster Version, Payload Mass (kg), and Launch site of 2015 launches where stage 1 failed to land on a drone ship.

There were two such occurrences.

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

---

```
%%sql
SELECT landing_outcome, COUNT(*) AS no_outcome
FROM SPACEXDATASET
WHERE landing_outcome LIKE 'Success%' AND DATE BETWEEN '2010-06-04' AND '2017-03-20'
GROUP BY landing_outcome
ORDER BY no_outcome DESC;
```

```
* ibm_db_sa://ftb12020:***@0c77d6f2-5da9-48a9-81f8-86b520b87518.bs2io90l08kqb1od8lcg
Done.
```

landing_outcome	no_outcome
Success (drone ship)	5
Success (ground pad)	3

This query returns a list of successful landings and between 2010-06-04 and 2017-03-20 inclusively.

There are two types of successful landing outcomes: drone ship and ground pad landings.

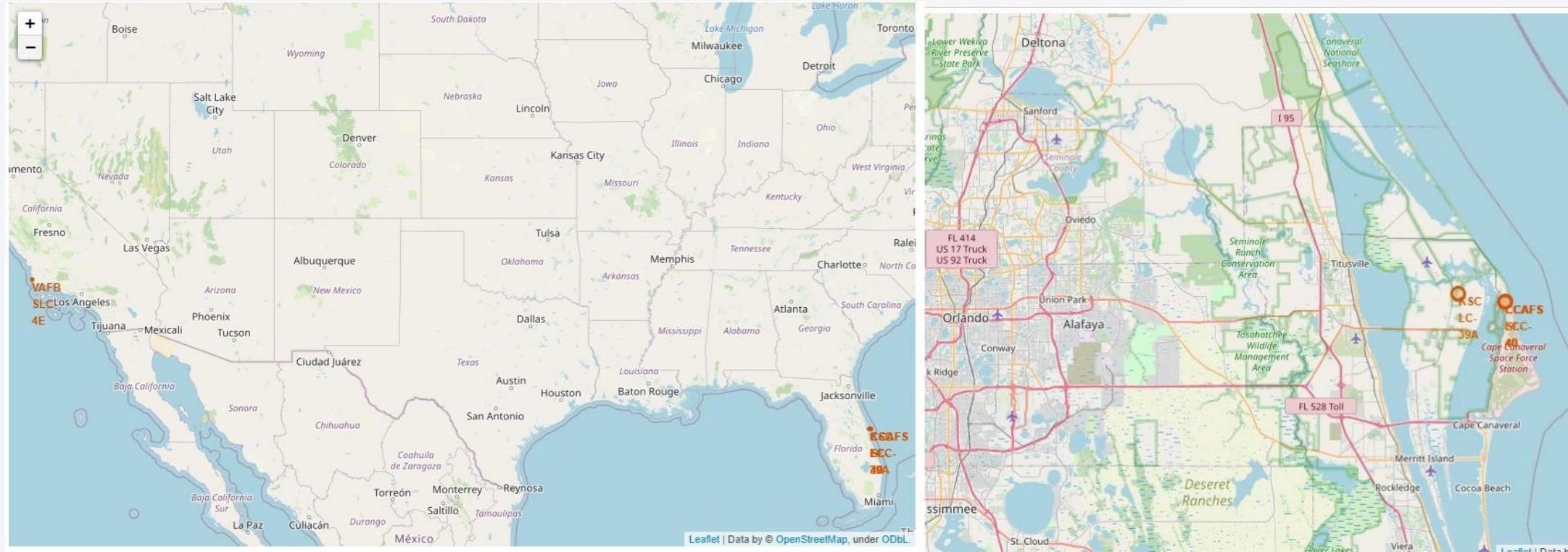
There were 8 successful landings in total during this time period

The background of the slide is a photograph taken from space at night. It shows the curvature of the Earth against a dark blue and black void of space. City lights are visible as small white dots and larger clusters of light, primarily concentrated in the lower right quadrant where the United States and Mexico would be. In the upper right, there is a bright green and yellow aurora borealis or aurora australis visible in the atmosphere.

Section 3

# Launch Sites Proximities Analysis

# Launch Site Locations



The left map shows all launch sites relative US map. The right map shows the two Florida launch sites since they are very close to each other. All launch sites are near the ocean.

# Launch Markers with Color-Coding

---

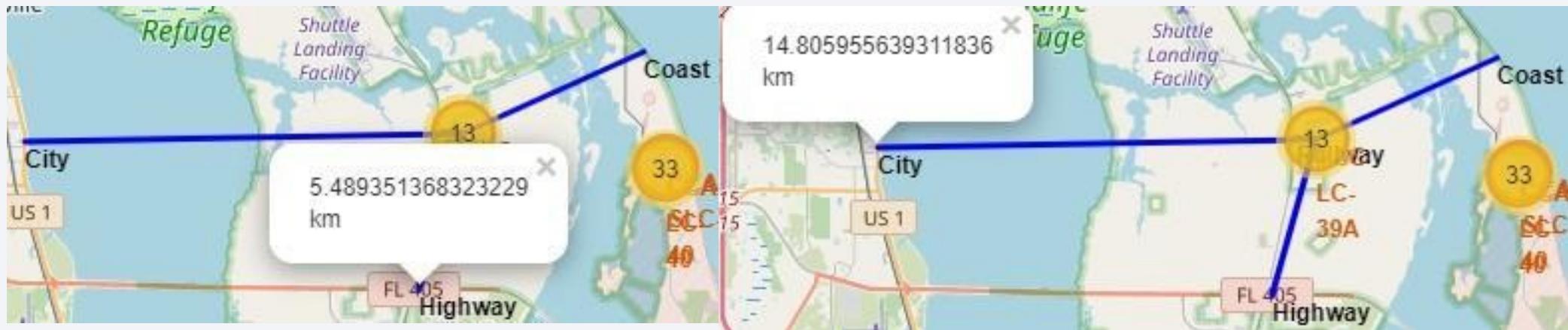


Clusters on Folium map can be clicked on to display each successful landing (green icon) and failed landing (red icon). In this example VAFB SLC-4E shows 4 successful landings and 6 failed landings.

# Proximities of Key Locations

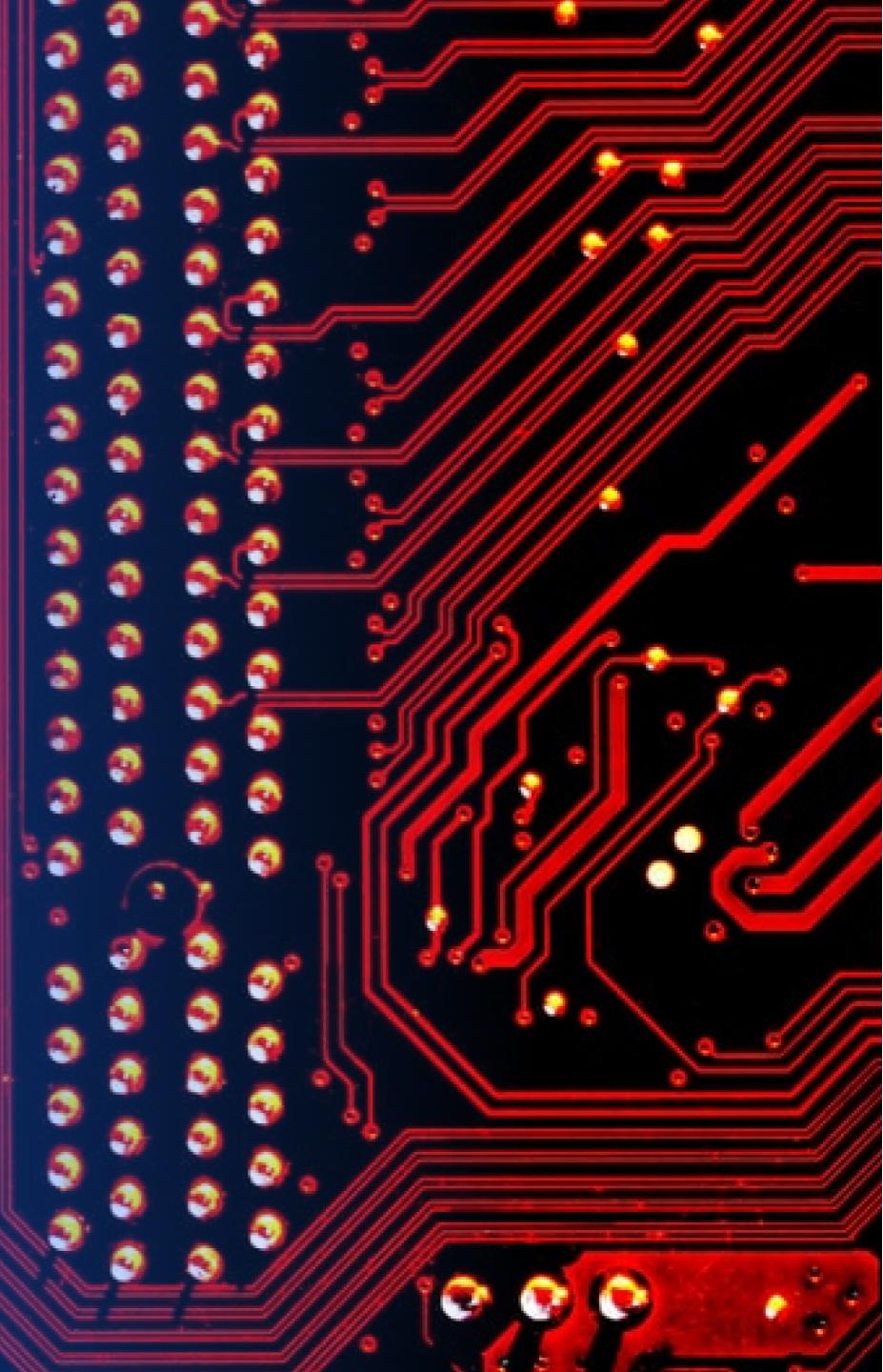


Using KSC LC-39A as an example, launch sites are very close to railways for large part and supply transportation. Launch sites are close to highways for human and supply transport. Launch sites are also close to coasts and relatively far from cities so that launch failures can land in the sea to avoid rockets falling on densely populated areas.



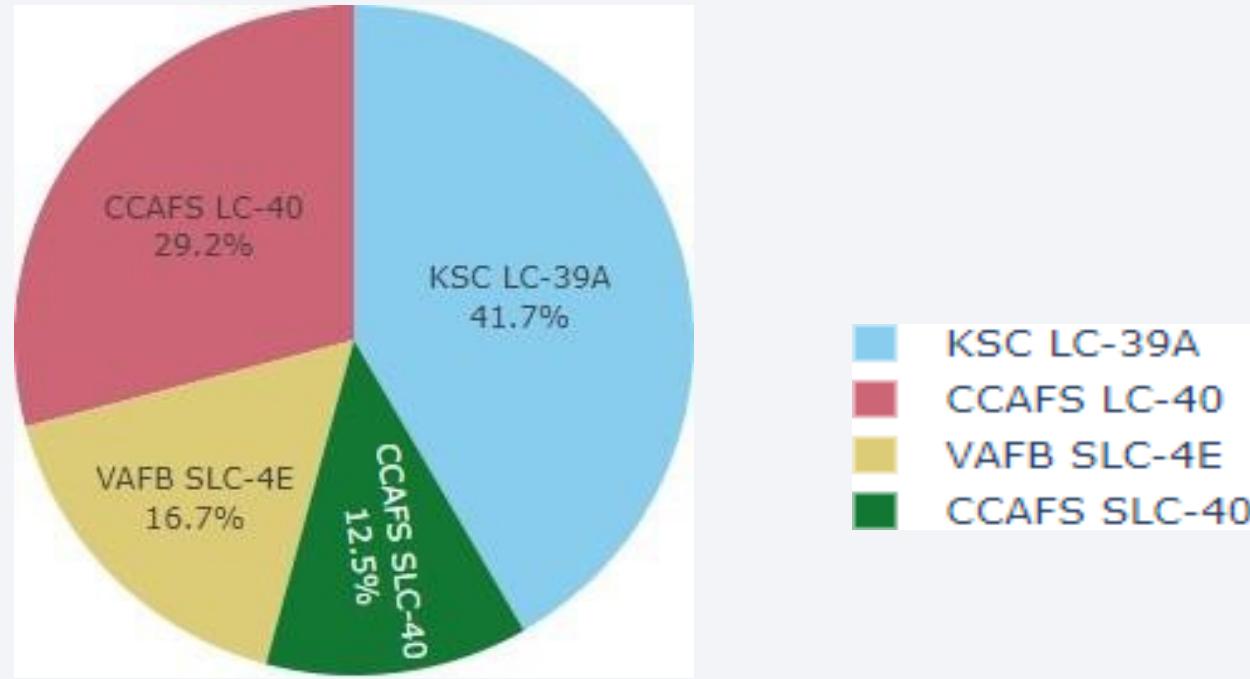
Section 4

# Build a Dashboard with Plotly Dash



# Successful Launches Across Launch Sites

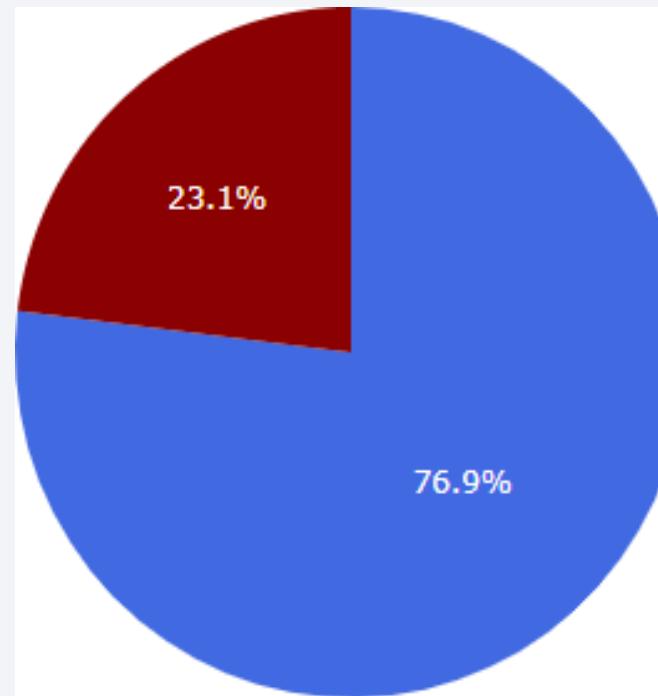
---



This is the distribution of successful landings across all launch sites. CCAFS LC-40 is the old name of CCAFS SLC-40 so CCAFS and KSC have the same amount of successful landings, but a majority of the successful landings were performed before the name change. VAFB has the smallest share of successful landings. This may be due to smaller sample and increase in difficulty of launching in the west coast.

# Highest Success Rate Launch Site

---



KSC LC-39A Success Rate (blue=success)



KSC LC-39A has the highest success rate with 10 successful landings and 3 failed landings.

# Payload Mass vs. Success vs. Booster Version Category



Plotly dashboard has a Payload range selector. However, this is set from 0-10000 instead of the max Payload of 15600. Class indicates 1 for successful landing and 0 for failure. Scatter plot also accounts for booster version category in color and number of launches in point size. In this particular range of 0-6000, interestingly there are two failed landings with payloads of zero kg.

Section 5

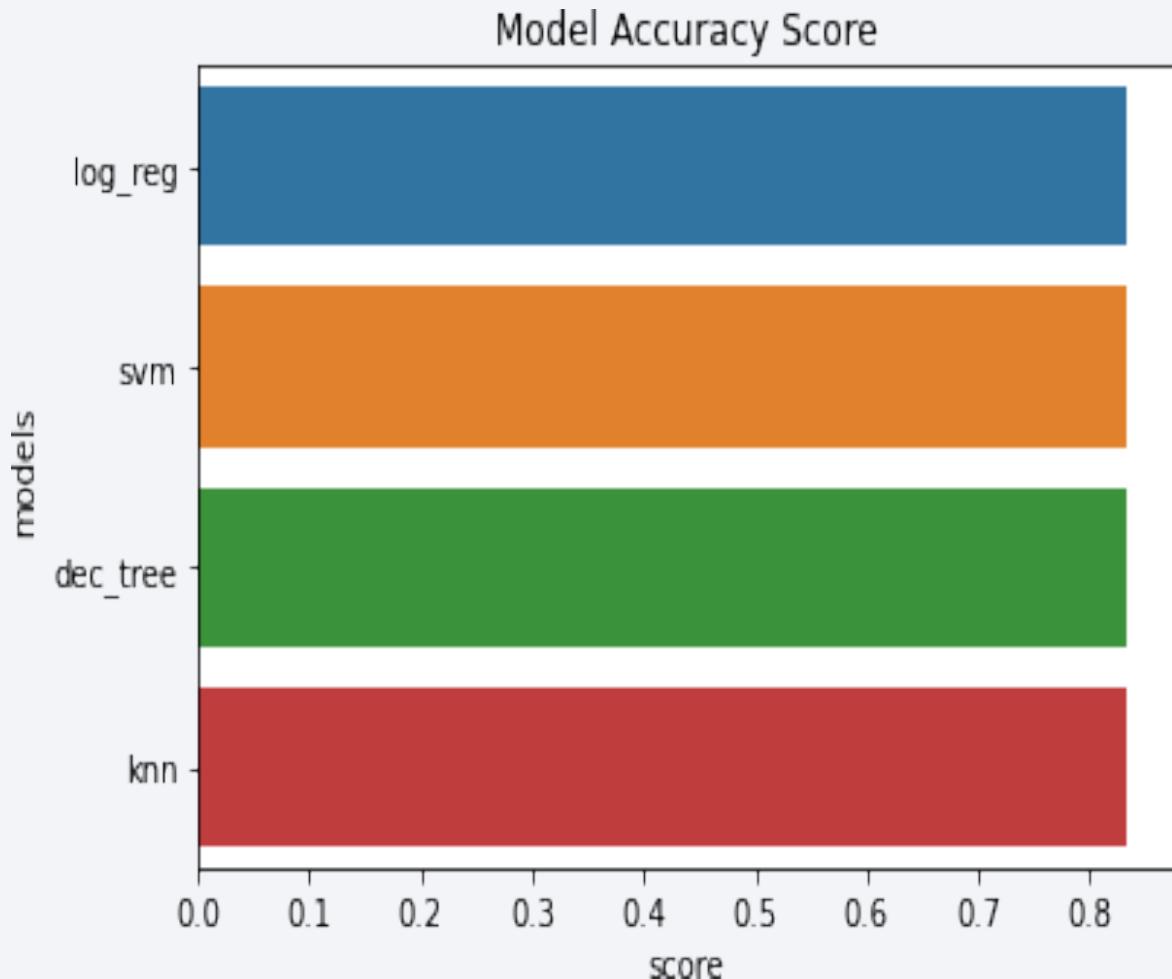
# Predictive Analysis (Classification)

# Classification Accuracy

---

All models had virtually the same accuracy on the test set at 83.33% accuracy. It should be noted that test size is small at only sample size of 18. This can cause large variance in accuracy results, such as those in Decision Tree Classifier model in repeated runs.

We likely need more data to determine the best model.



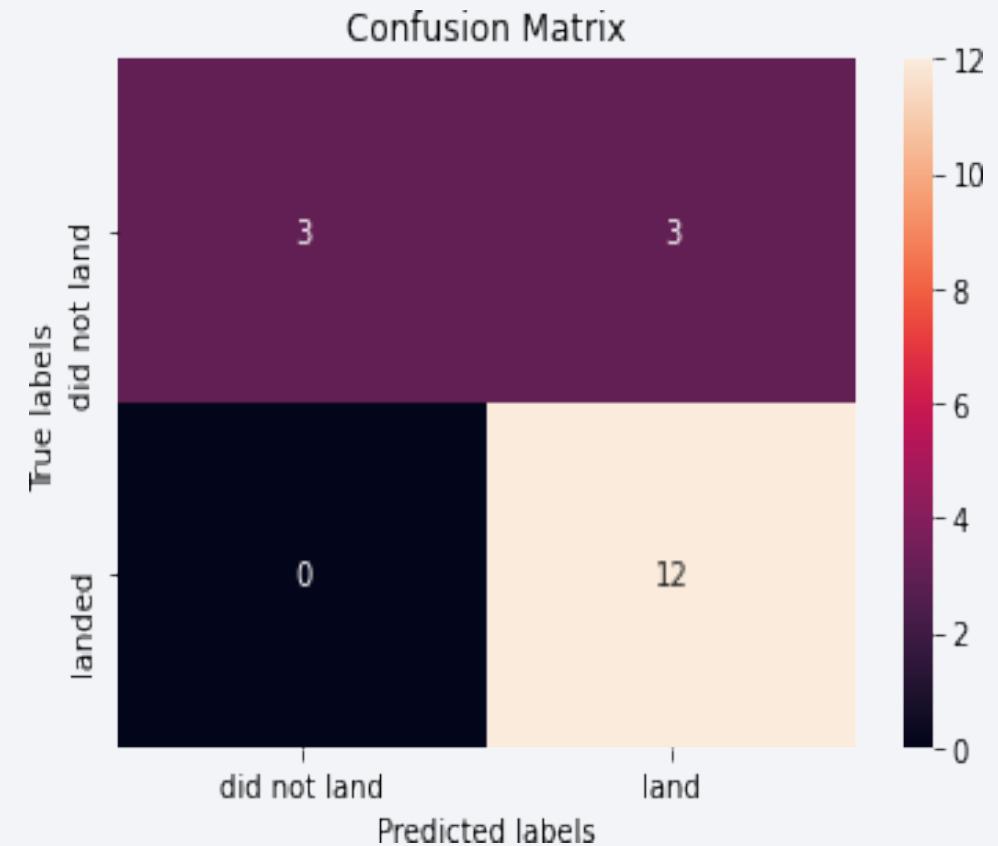
# Confusion Matrix

---

Since all models performed the same for the test set, the confusion matrix is the same across all models. The models predicted 12 successful landings when the true label was successful landing.

The models predicted 3 unsuccessful landings when the true label was unsuccessful landing.

The models predicted 3 successful landings when the true label was unsuccessful landings (false positives). Our models over predict successful landings.



Correct predictions are on a diagonal from top left to bottom right.

# Conclusions

---

- Objective - develop a machine learning model for Space Y, aiming to compete with SpaceX.
- Model's aim - forecast the successful landing of Stage 1 to potentially save approximately \$100 million USD.
- To achieve this following steps were taken,
  - Data collection - sources from SpaceX API and web scraping of SpaceX Wikipedia.
  - Data preparation - labeling and storing data into a DB2 SQL database.
  - Data visualization - plotting and making dashboard for data visualization.

## ● **Machine Learning Model results an accuracy rate of 83.33%.**

- This model enables Allon Mask of SpaceY to predict, with considerable accuracy, whether a launch will feature a successful Stage 1 landing prior to its execution.
- This prediction aids in deciding whether the launch should proceed or not.

Furthermore,

It has been always recommended acquiring additional data to fine-tune the machine learning model and boost its predictive accuracy.

# Appendix

---

## Credit and Acknowledgments:

Project Lead: [Rav Ahuja](#)

Instructional Designer: [Lakshmi Holla](#)

Lab Authors: [Joseph Santarcangelo](#), [Yan Luo](#), [Azim Hirjani](#), [Lakshmi Holla](#)

Technical Advisor: [Yan Luo](#)

Publishing: [Grace Barker](#), [Rachael Jones](#)

Project Coordinators: [Kathleen Bergner](#)

Narration: [Bella West](#)

Video Production: [Simer Preet](#), [Lauren Hall](#), [Hunter Bay](#), [Tanya Singh](#), [Om Singh](#)

Teaching Assistants and Forum Moderators: [Malika Singla](#), [Duvvana Mrutyunjaya Naidu](#), [Lakshmi Holla](#), [Anita Verma](#)

Special Thanks to Coursera for this opportunity.

Github URL: [https://github.com/MaryamBeik/testrepo/blob/main/SpaceX\\_Machine%20Learning%20Prediction\\_Part\\_5.ipynb](https://github.com/MaryamBeik/testrepo/blob/main/SpaceX_Machine%20Learning%20Prediction_Part_5.ipynb)

Thank you!

