

# Winning Space Race with Data Science

Alin Rizea

<https://github.com/alin-r-13/IBMDatascienceProject>

24.11.2021



# Outline

---

- Executive Summary – slide 3
- Introduction – slide 4
- Methodology – slide 6
- Results – slide 16
- Conclusion – slide 45
- Appendix – slide 46

# Executive Summary

---

- Collected data from public SpaceX API and SpaceX Wikipedia page. Created labels column “class” which classifies successful landings. Explored data using SQL, visualization, Folium maps and dashboards. Gathered relevant columns to be used as features. Changed all categorical variables to binary using one hot encoding. Standardized data and used GridSearchCV to find the best parameters for machine learning models. Visualized accuracy score of all models.
- Four machine learning models were produced: Logistic Regression, Support Vector Machines, Decision Tree Classifier and K Nearest Neighbors. All models produced similar results with an accuracy rate of 83.33%. Overall, pretty much all models over predicted successful landings. Therefore, more data is needed for better model determination and accuracy.

# Introduction

---



SpaceX Falcon 9 Rocket – The Verge

## Background:

- Commercial Space Age is here
- Space X has best pricing (\$62 mil. Vs. \$165 mil)
- Largely due to the ability to recover part of the rocket (Stage 1)
- Space Y wants to compete with Space X

## Problem:

- Space Y tasks us to train a machine learning model in order to predict successful Stage 1 recovery.

Section 1

# Methodology

# Methodology

---

## Executive Summary

- Data collection methodology:
  - Combined data from SpaceX public API and SpaceX Wikipedia page.
- Perform data wrangling
  - Data was processed by classifying true landings as successful and unsuccessful otherwise.
- 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 models using GridSearchCV

# Data Collection

---

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

The next slide will show the flowchart of data collection from API and the one after will show the flowchart of data collection from webscraping.

## SpaceX API Data Columns:

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

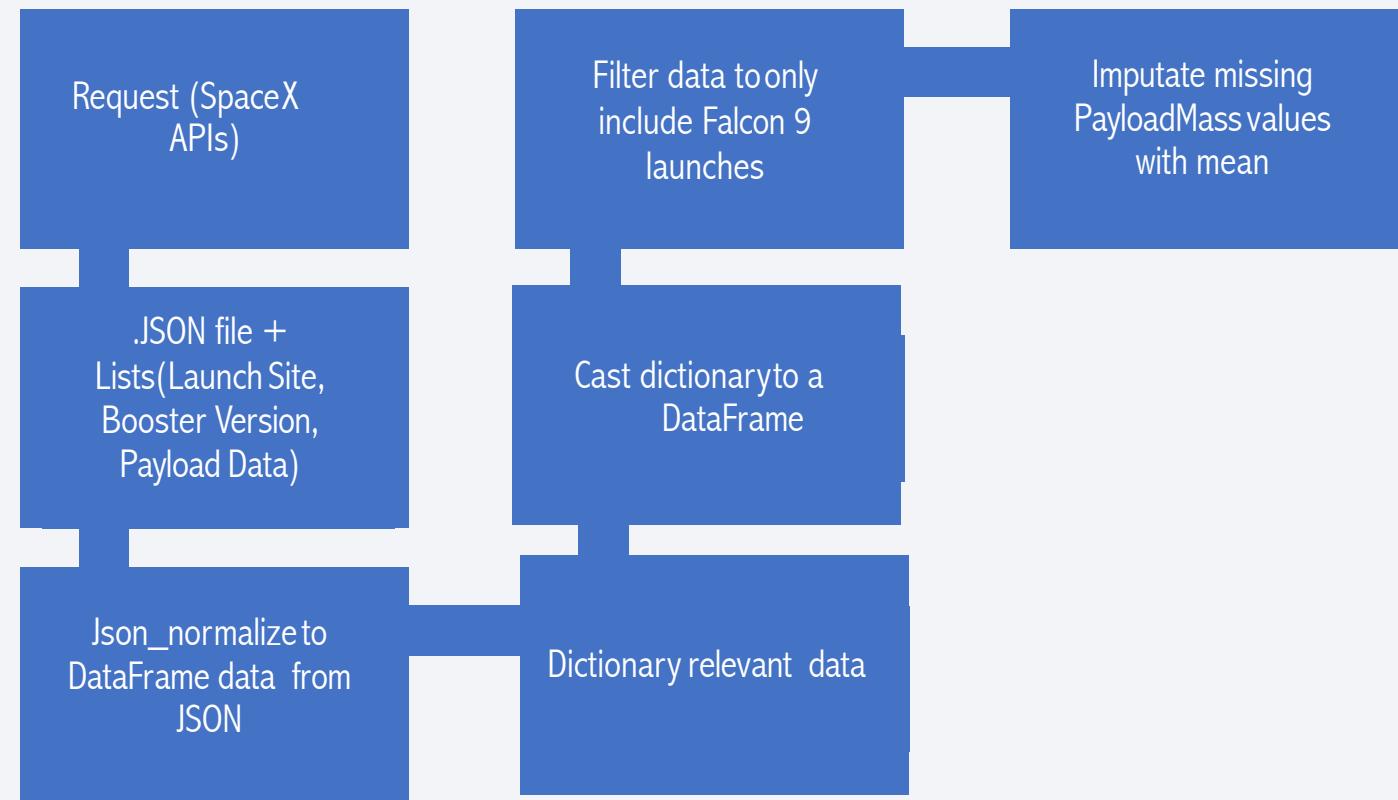
## Wikipedia Webscraping Data Columns:

Flight No., Launch Site, Payload, PayloadMass, Orbit, Customer, Launch outcome, Version Booster, Booster Landing, Date, Time

# Data Collection – SpaceX API

## GitHub url:

<https://github.com/alin-r-13/IBMDatascienceProject/blob/master/10.%20%20Data%20Science%20and%20Machine%20Learning%20Capstone%20Project/Week%201%20Introduction/Data%20Collection%20API%20.ipynb>

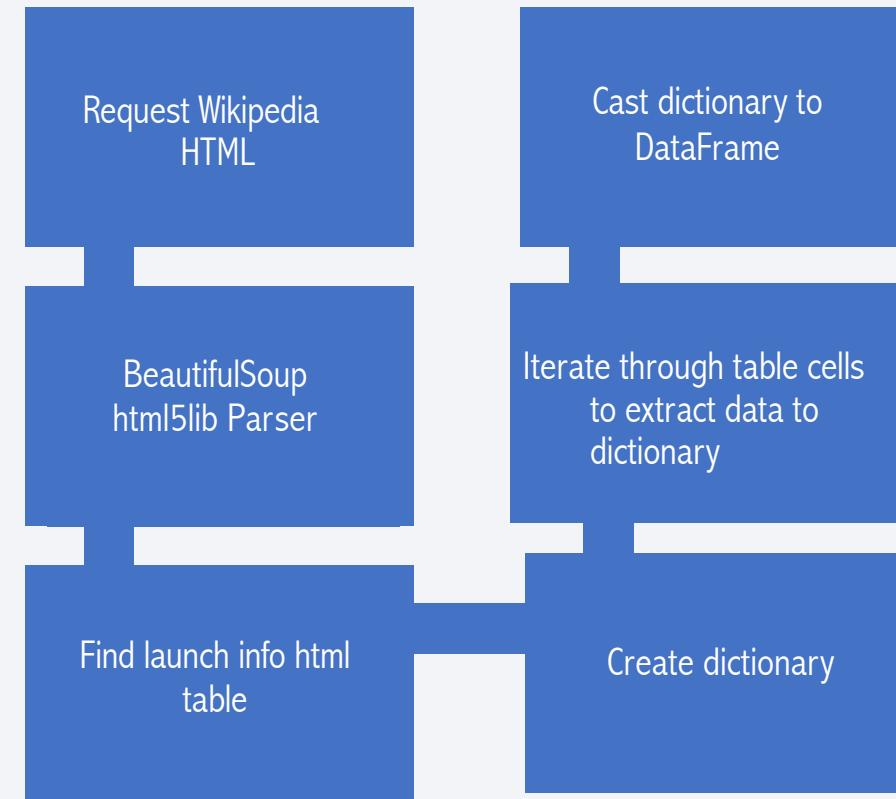


# Data Collection - Scraping

---

## GitHub url:

<https://github.com/alin-r-13/IBMDatascienceProject/blob/master/10.%20%20Data%20Science%20and%20Machine%20Learning%20Capstone%20Project/Week%201%20Introduction/Data%20Collection%20with%20Web%20Scraping.ipynb>



# Data Wrangling

---

## Flow explained

Created a training label 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/alin-r-13/IBMDatascienceProject/blob/master/10.%20%20Data%20Science%20and%20Machine%20Learning%20Capstone%20Project/Week%201%20Introduction/Data%20wrangling%20.ipynb>

# EDA with Data Visualization

---

Exploratory Data Analysis performed on variables Flight Number, Payload Mass, Launch Site, Orbit, Class and Year.

## **Plots used:**

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

Scatter plots, line charts and bar plots were used to compare relationships between variables in order to decide if a relationship exists so that they could be used in training the machine learning model.

## **GitHub url:**

<https://github.com/alin-r-13/IBMDatascienceProject/blob/master/10.%20%20Data%20Science%20and%20Machine%20Learning%20Capstone%20Project/Week%202%20EDA/EDA%20with%20Visualization.ipynb>

# EDA with SQL

---

- Loaded dataset 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/alin-r-13/IBMDatascienceProject/blob/master/10.%20%20Data%20Science%20and%20Machine%20Learning%20Capstone%20Project/Week%202%20EDA%20with%20SQL.ipynb>

# Build an Interactive Map with Folium

---

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/alin-r-13/IBMDatascienceProject/blob/master/10.%20%20Data%20Science%20and%20Machine%20Learning%20Capstone%20Project/Week%203%20Interactive%20Visual%20Analytics%20and%20Dashboard/Interactive%20Visual%20Analytics%20with%20Folium.ipynb>

# Build a Dashboard with Plotly Dash

---

Dashboard includes a pie chart and a scatter plot.

Pie chart can be selected to show distribution of successful landings across all launch sites and can be selected to show individual launch site success rates.

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

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.

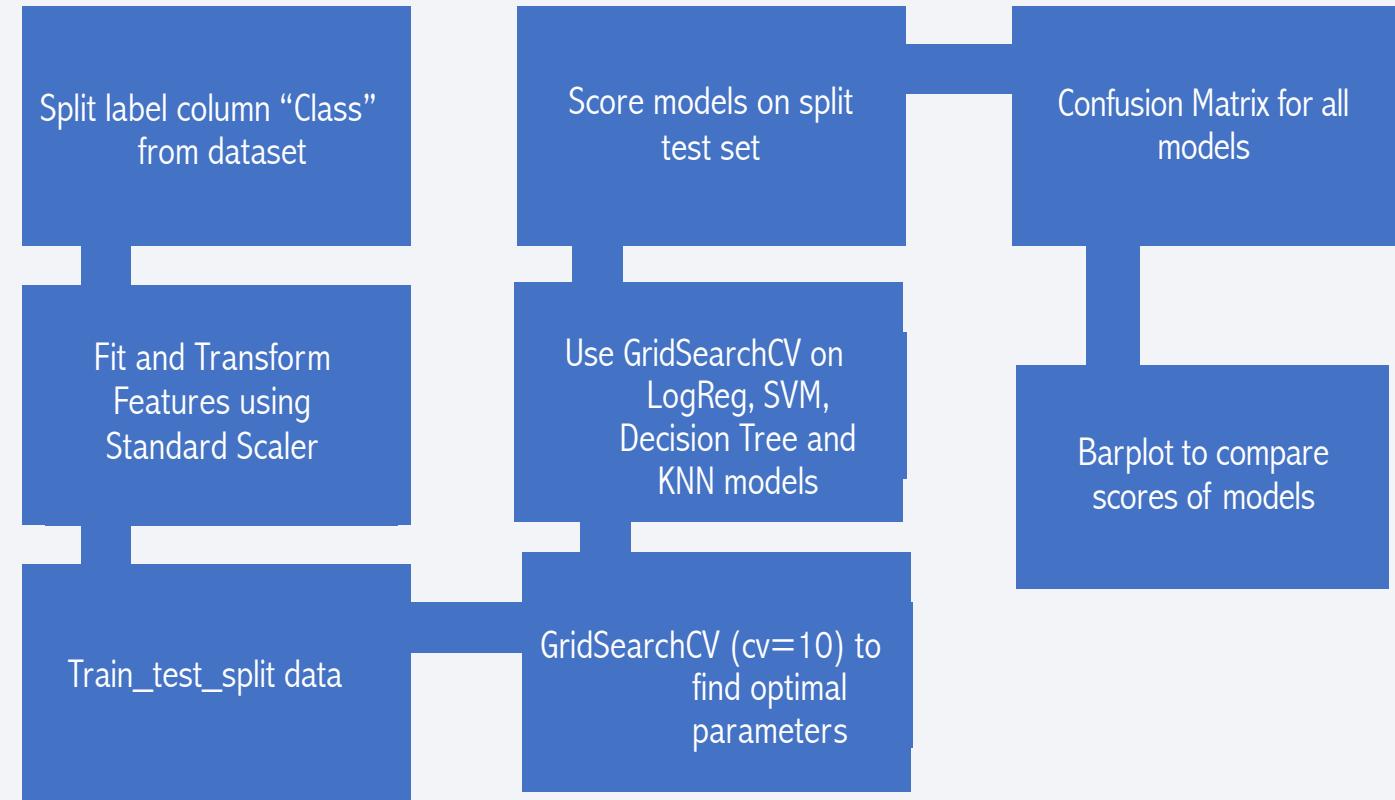
## **GitHub url:**

[https://github.com/alin-r-13/IBMDatascienceProject/blob/master/10.20%20Data%20Science%20and%20Machine%20Learning%20Capstone%20Project/Week%203%20Interactive%20Visual%20Analytics%20and%20Dashboard/spacex\\_dash\\_app.py](https://github.com/alin-r-13/IBMDatascienceProject/blob/master/10.20%20Data%20Science%20and%20Machine%20Learning%20Capstone%20Project/Week%203%20Interactive%20Visual%20Analytics%20and%20Dashboard/spacex_dash_app.py)

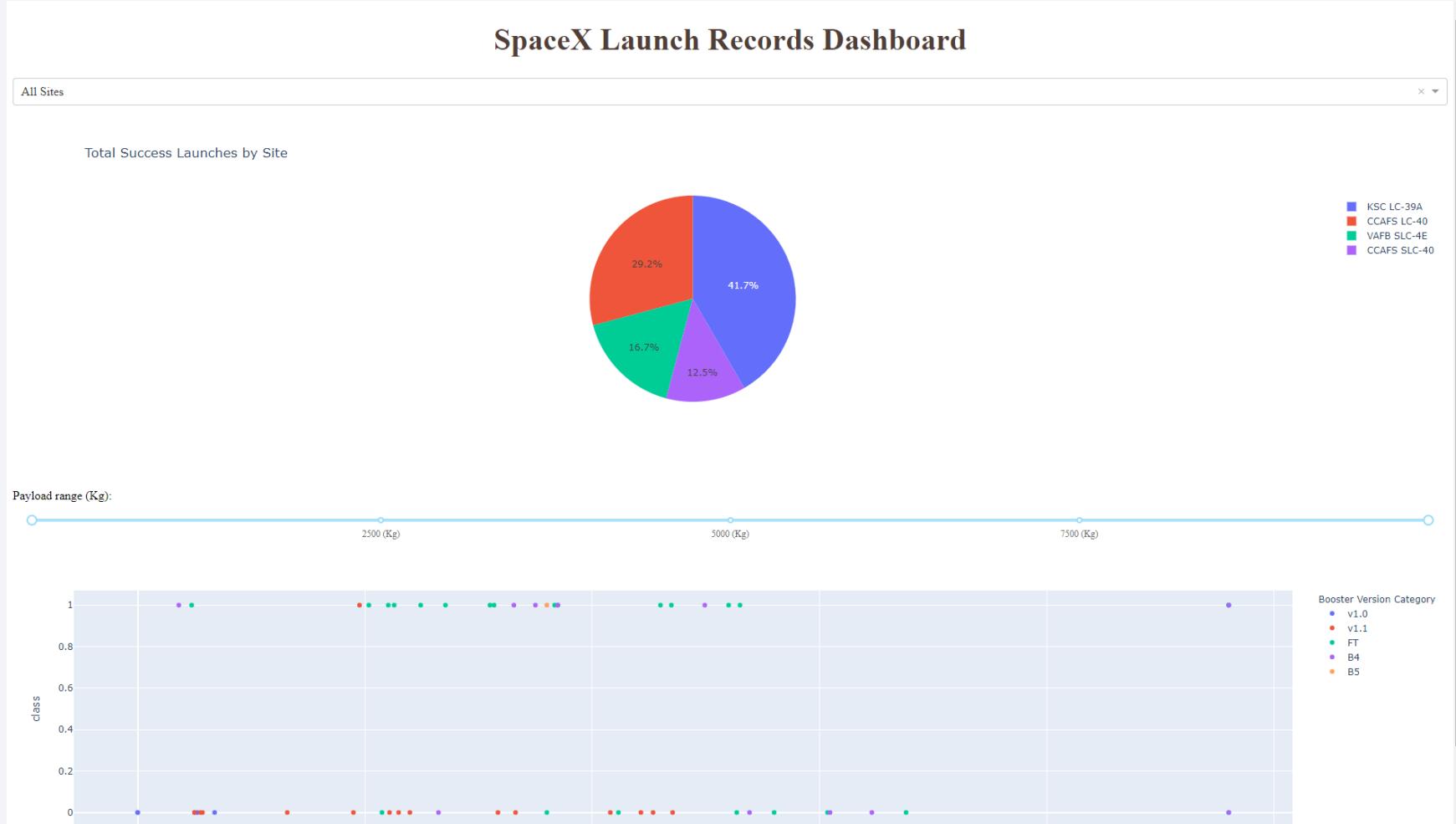
# Predictive Analysis (Classification)

## GitHub url:

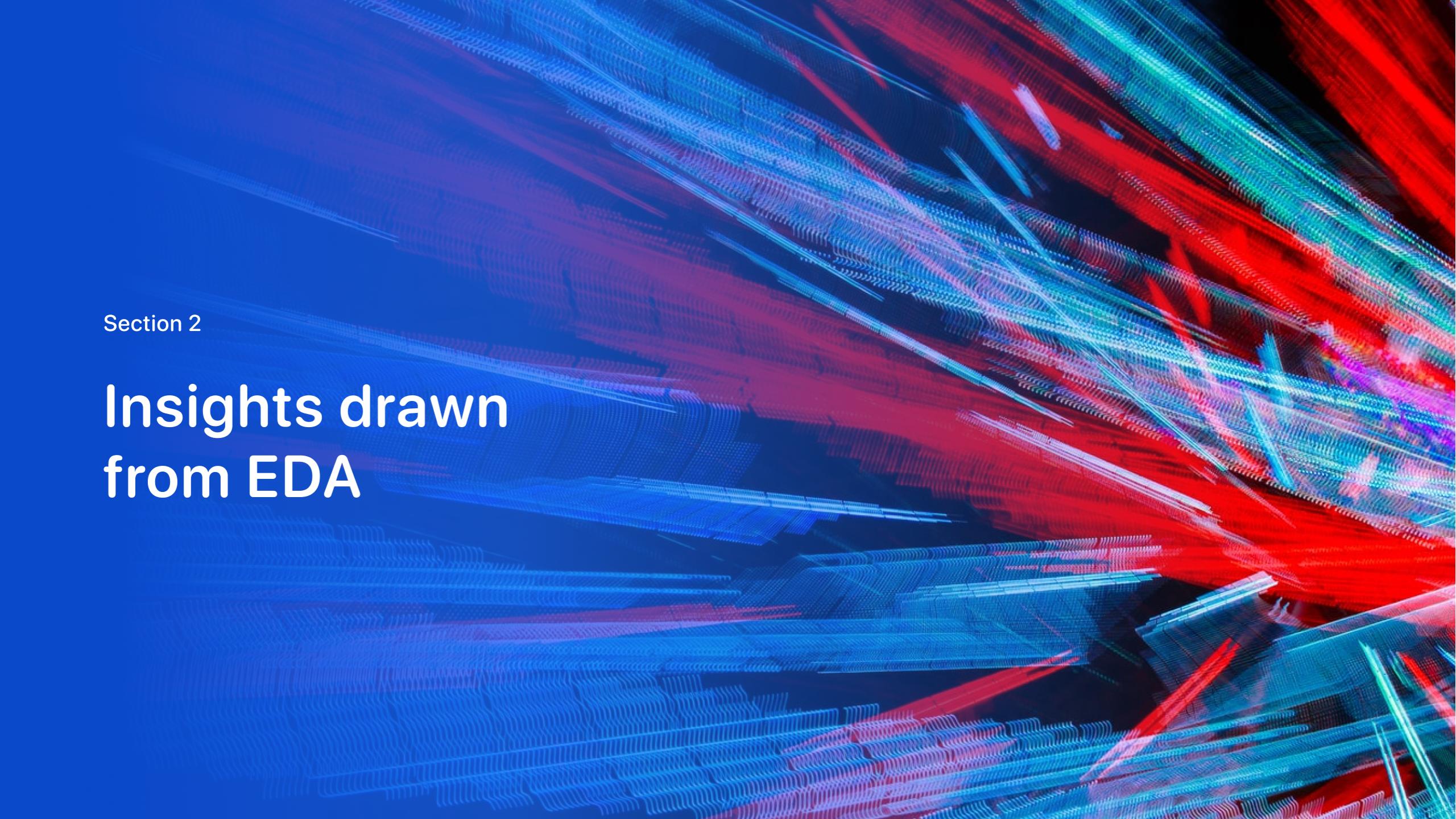
[https://github.com/alin-r-13/IBMDatascienceProject/blob/master/10.%20%20Data%20Science%20and%20Machine%20Learning%20Capstone%20Project/Week%204%20Predictive%20Analysis%20\(Classification\)/Machine%20Learning%20Prediction.ipynb](https://github.com/alin-r-13/IBMDatascienceProject/blob/master/10.%20%20Data%20Science%20and%20Machine%20Learning%20Capstone%20Project/Week%204%20Predictive%20Analysis%20(Classification)/Machine%20Learning%20Prediction.ipynb)



# Results



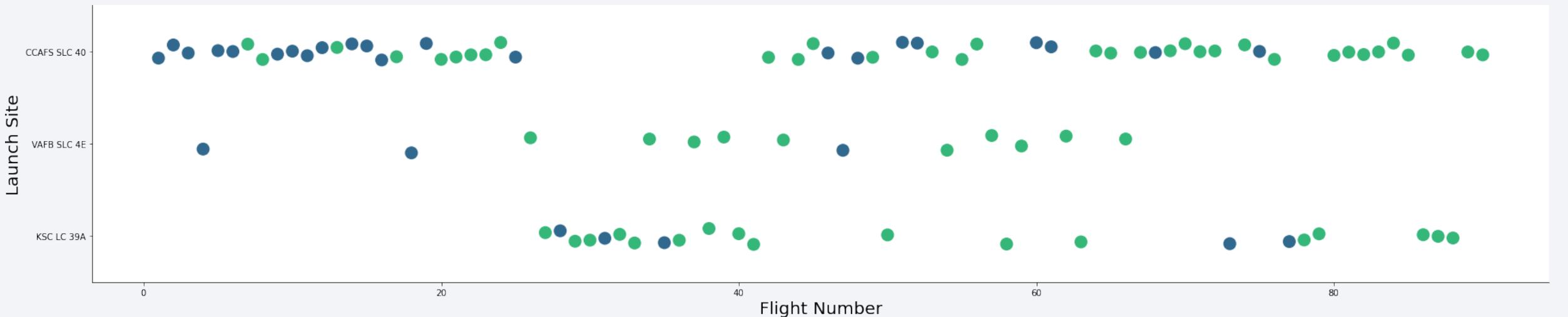
This is a preview of the Plotly dashboard. The next slides will show the results of EDA with visualization, EDA with SQL, Interactive Map with Folium and finally, the results of our model with around 83% accuracy.

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



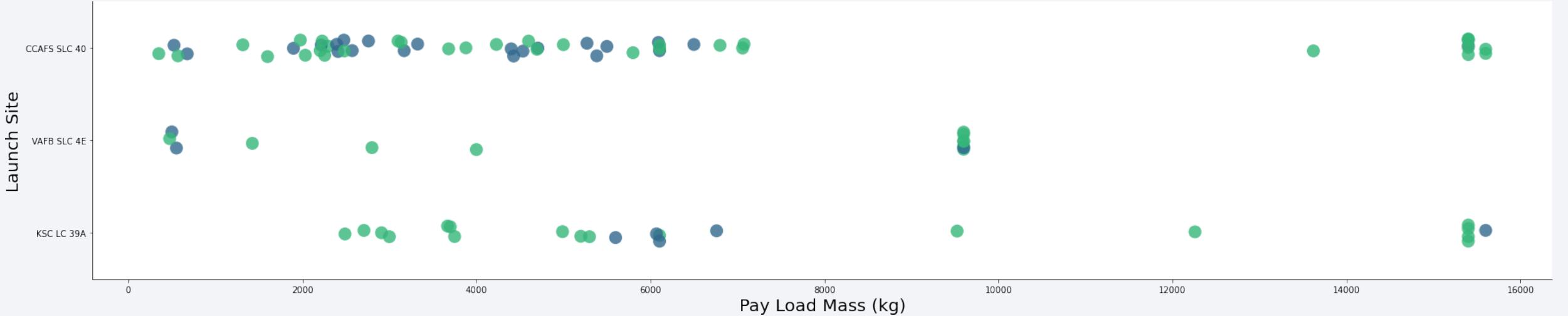
Green indicates successful launch; Purple indicates unsuccessful launch.

Graphic suggests an increase in success rate over time (indicated in Flight Number).

Likely a great breakthrough around flight 20 which significantly increased the success rate.

CCAFS seems to be the main launch site as it has the most volume.

# Payload vs Launch Site

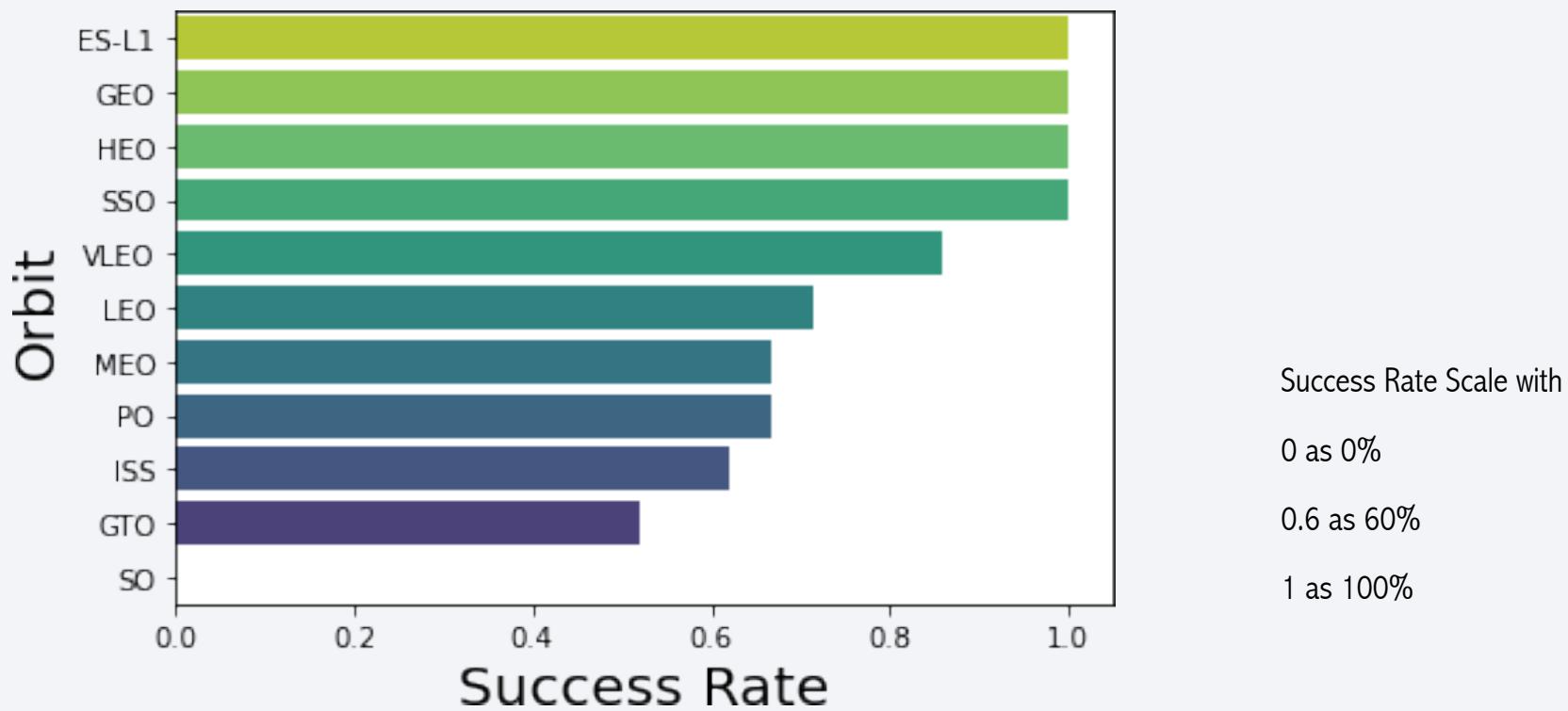


Green indicates successful launch; Purple indicates unsuccessful launch.

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)

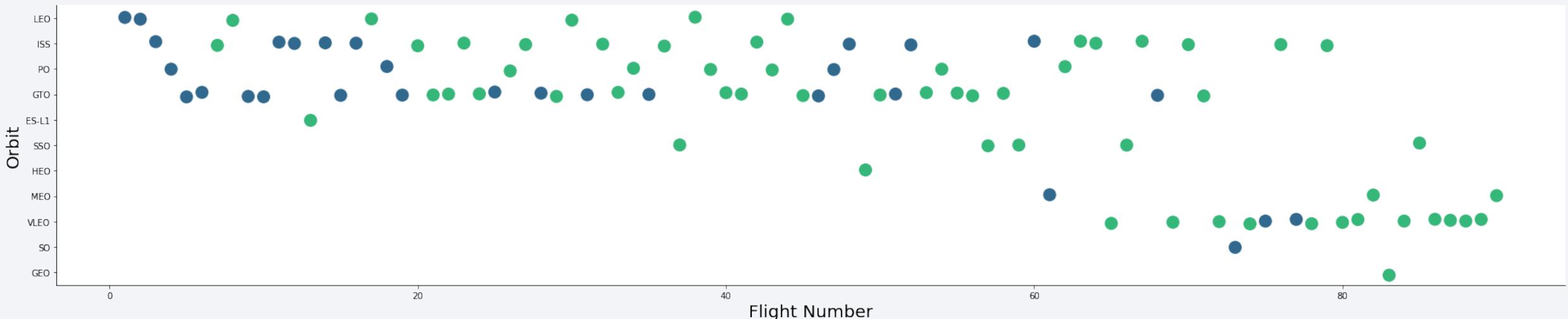
SSO(5) has 100% success rate

VLEO(14) has decent success rate and attempts

SO(1) has 0% success rate

GTO(27) has around 50% success rate, but larges sample

# 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 and then switched 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.

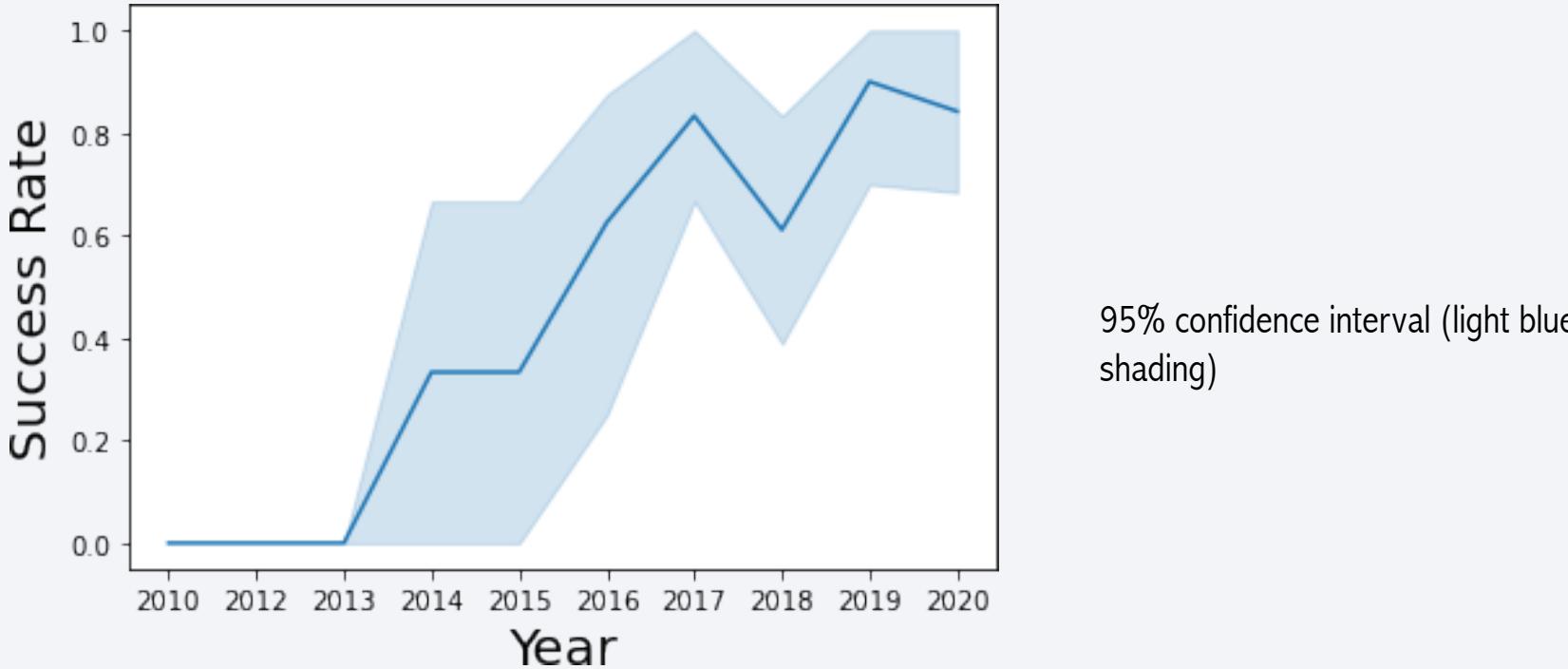
Payload Mass seems to correlate with orbit.

HEO 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 decrease in 2018.

Success in recent years is around 80%.

# All Launch Site Names

---

In [4]:

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

Out[4]:

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

- Queried unique launch site named from the database.
- CCAFS SLC-40 and CCAFSSLC-40 are likely representing 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'

In [5]:

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

\* ibm\_db\_sa://ftb12020:\*\*\*@0c77d6f2-5da9-48a9-81f8-86b520b87518.bs2io90l08kqb1od81cg.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

First five entries in database with Launch Site name starting with 'CCA'.

# Total Payload Mass 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

---

```
%%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 or launches which used booster version F9 v1.1.

Average payload mass of F9 v1.1 is on the low end of our payload mass range.

# First Successful Ground Landing Date

```
%%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

The 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 with 2014.

# Successful Drone Ship Landing with Payload between 4000 and 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.bs2io90l08kqb1od8lcg.database.
Done.
```

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

The 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

```
%%sql
SELECT mission_outcome, COUNT(*) AS no_outcome
FROM SPACEXDATASET
GROUP BY mission_outcome;
```

```
* ibm_db_sa://ftb12020:***@0c77d6f2-5da9-48a9-1
Done.
```

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

The 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 That Carried Maximum Payload

```
%%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 that the payload mass correlates with the booster version that is used.

# 2015 Failed Drone Ship Landing Records

---

```
%%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.bs2io90108kqb1od81cg.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 a drone ship.

There were two occurrences.

## Ranking Counts of Successful Landings 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.bs2io90108kqb1od81ce  
Done.
```

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

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

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

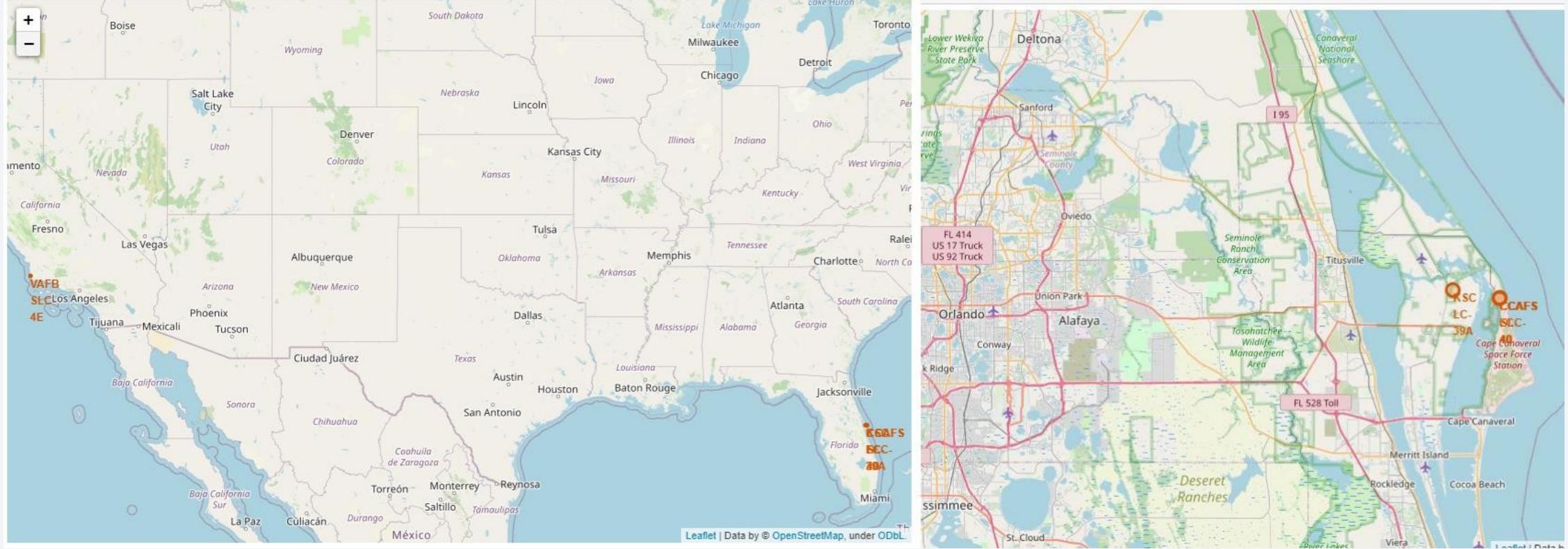
There were 8 successful landings in total for the time period mentioned above.

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 appears. In the upper right, there are bright green and yellow bands of the Aurora Borealis (Northern Lights) dancing across the sky.

Section 4

# 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.

# Color-Coded Launch Markers

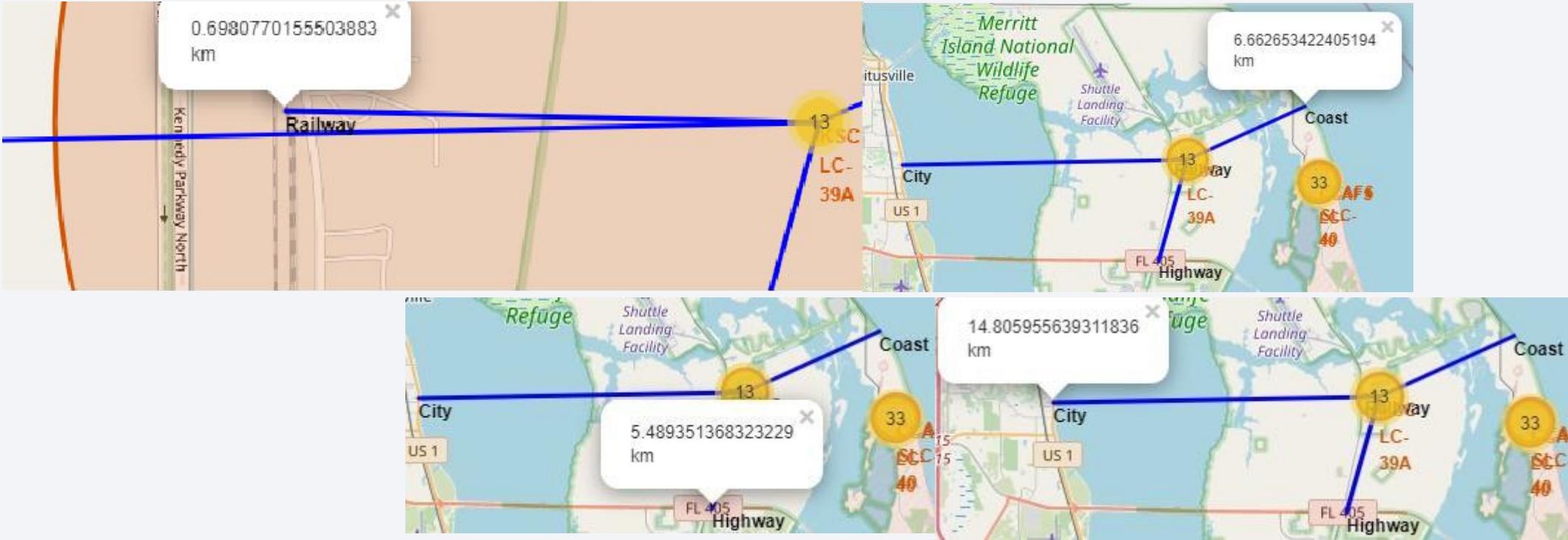
---



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.

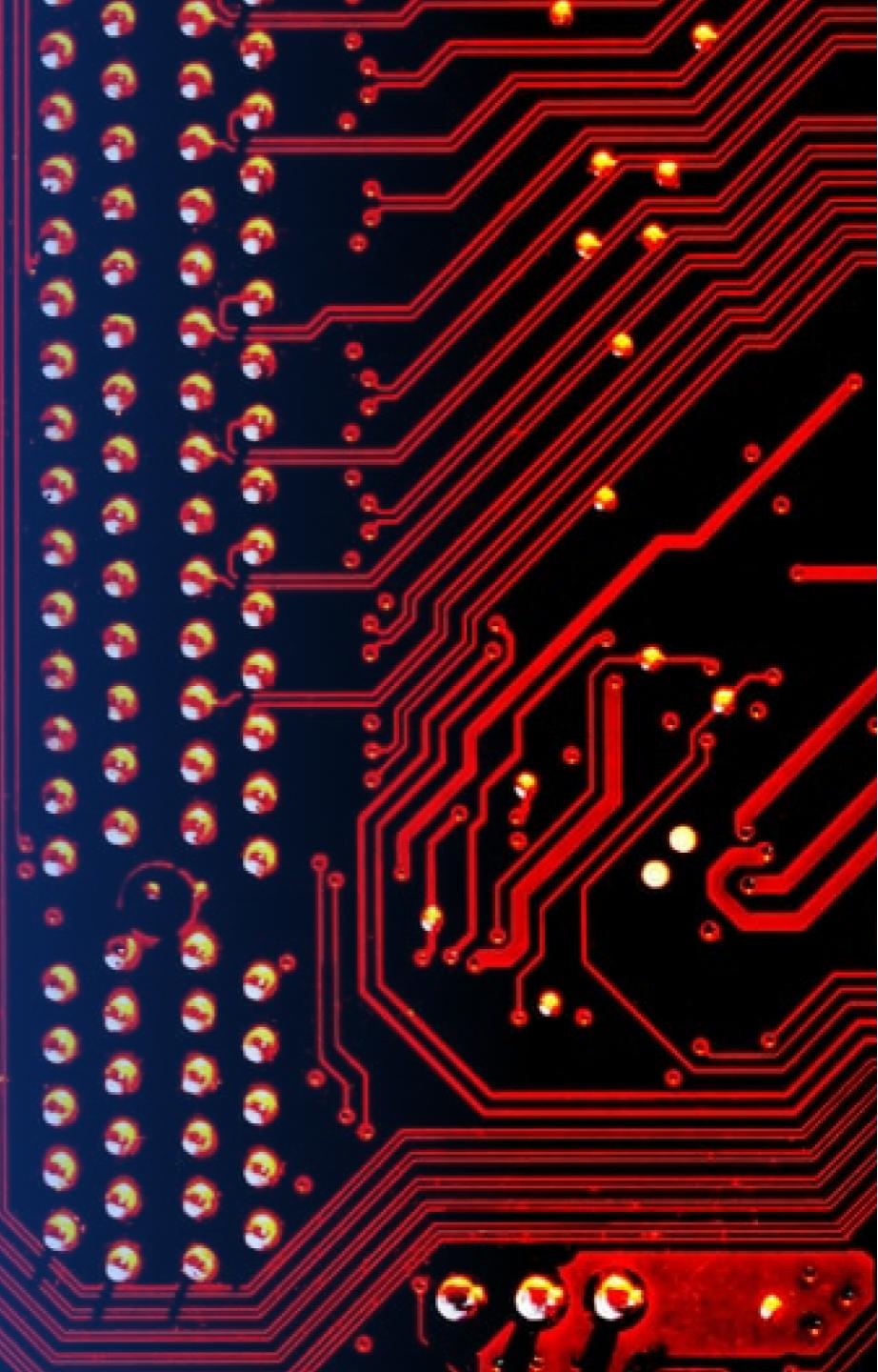
# Key Location Proximities



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 the coasts and relatively far from the cities so that launch failures can land in the sea to avoid rockets falling on densely populated areas.

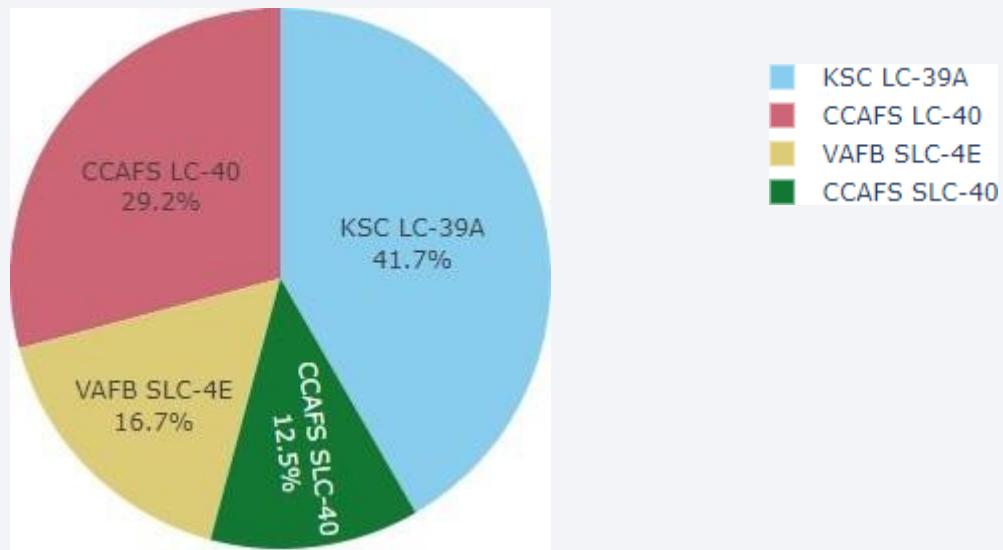
Section 5

# Build a Dashboard with Plotly Dash



# Successful Launches Across Launch Sites

---



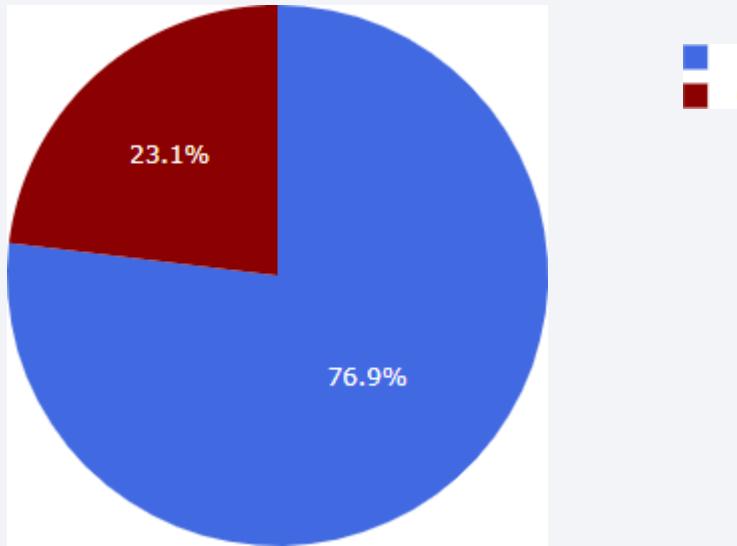
This is the distribution of successful landings across all launch sites. CCAFS LC-40 is the old name of CCAFS SLC-40, so CCAFS and KSC have the same amount of successful landings, but a majority of the successful landings were performed before the name change.

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

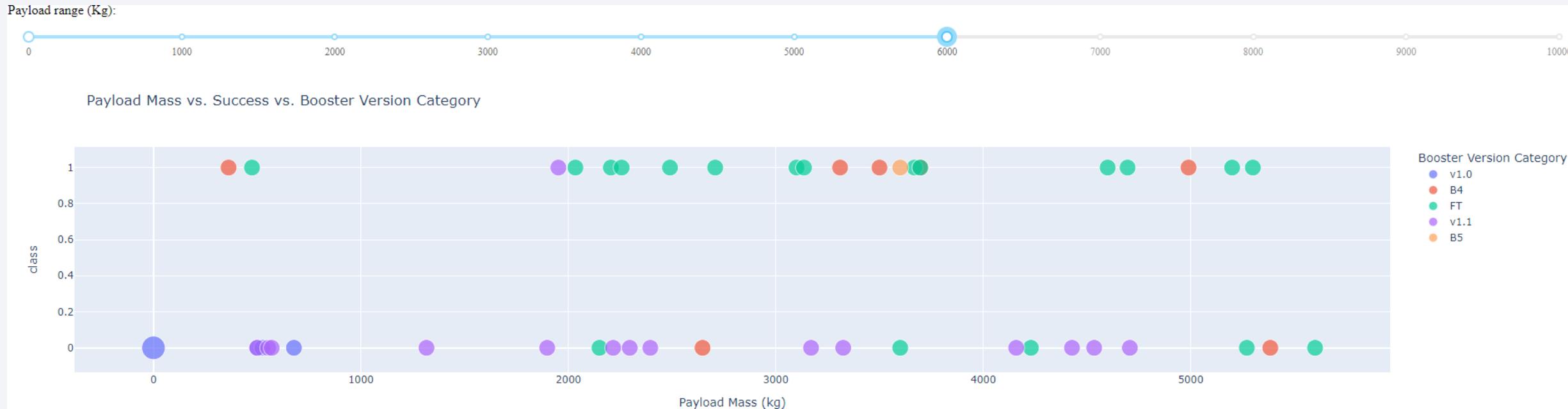
---

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. The scatter plot also accounts for booster version category in color and number of launches in point size.

In this particular range of 0-6000, it seems there are two failed landings with payloads of zero kg.

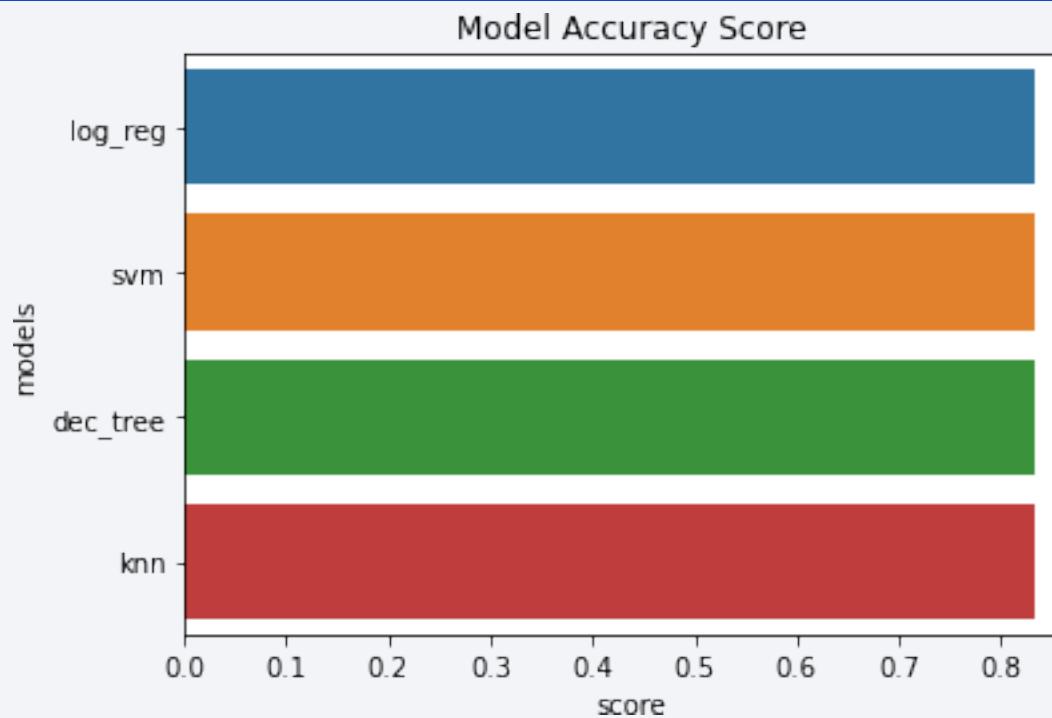
The background of the slide features a dynamic, abstract design. It consists of several thick, curved lines that transition from a bright yellow at the top right to a deep blue at the bottom left. These lines create a sense of motion and depth, resembling a tunnel or a stylized landscape. The overall effect is modern and professional.

Section 6

# Predictive Analysis (Classification)

# Classification Accuracy

---



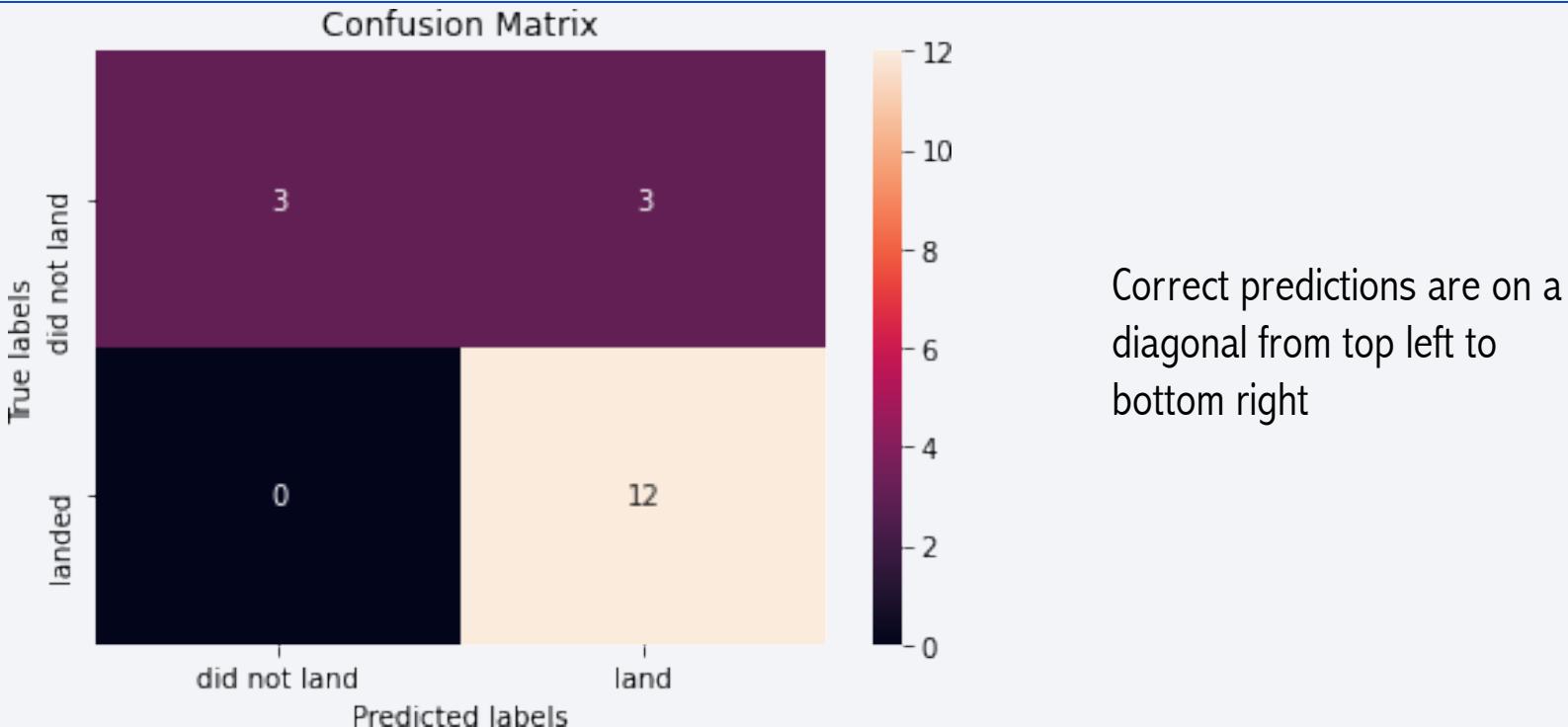
All models had around the same accuracy on the test set, which was at 83.33% accuracy.

However, it should be noted that the test size is small, with a sample size of 18.

This can cause large variance in accuracy results, such as those in Decision Tree Classifier model in repeated runs.

It is very likely that we will 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 success 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.

# Conclusions

---

- Our task was to develop a machine learning model for Space Y who wants to bid against SpaceX.
- The goal of the model is to predict when Stage 1 will successfully land to save ~\$100 million.
- So, for that we used data from a public SpaceX API and web scraping SpaceX Wikipedia page.
- We created data labels and stored data into a DB2 SQL database.
- We created a dashboard for visualization.
- We created a machine learning model with an accuracy of 83%
- Allon Mask of Space Y can use this model to predict with relatively high accuracy whether a launch will have a successful Stage 1 landing before launch to determine whether the launch should be made or not.
- If possible, more data should be collected to better determine the best machine learning model and improve accuracy.

# Appendix

---

**GitHub url:**

<https://github.com/alin-r-13/IBMDatascienceProject>

**Instructors:**

Romeo Kienzler, Saeed Aghabozorgi, Rav Ahuja, Joseph Santarcangelo, Saishruthi Swaminathan, Azim Hirjani, Svetlana Levitan, Maureen McElaney, Alex Akison, Yan Luo

Thank you!

