

# Winning Space Race with Data Science

K.Santhana Gopala  
Krishnan  
23.01.23



# Outline

---

- Executive Summary
- Introduction
- Methodology
- Results
- Conclusion
- Appendix

# Executive Summary

---

- **Summary of methodologies**

## **Data Collection using API (SpaceX) and Web Scraping**

- **Data Wrangling (fill Falcon 9 the missing data, include successful launches for reusing)**
- **Exploratory Data Analysis (EDA) with SQL and Data Visualization (identify features that have influence over the success of launch)**
- **Visual Analytics through Plotly Dash and Folium (look at the launch sites on the map relevant to its proximities)**
- **Machine Learning Prediction (with 4 classification)**
- **Summary of all results**
- **Exploratory Data Analysis result**
- **Interactive analytics in screenshots**
- **Predictive Analytics result**

# Introduction

---

- Project background and context
  - Space X 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 Space X can reuse the first stage. Therefore, if we can determine that the first stage will land, we can define the cost of a launch. This information can be used if other companies want to attempt against space X for a rocket launch. This goal of the project is to create a machine learning pipeline to predict if the first stage will land successfully
- Problems you want to find answers
  - - What factors determine if the rocket will land successfully?
  - - determine the price of each launch by gathering information about Space X and creating dashboards

Section 1

# Methodology

# Methodology

---

## Executive Summary

- Data collection methodology:
  - Space X API (<https://api.spacexdata.com/v4/rockets/>)
  - WebScraping (Wikipedia)
- Perform data wrangling
- data was categorized by a label (landing outcome) based on outcome data after data wrangling and analyzing attributes
- Perform exploratory data analysis (EDA) using visualization and SQL
- Perform interactive visual analytics using Folium and Plotly Dash
- Perform predictive analysis using classification models
  - Data were normalized, split into train and test data sets then evaluated by different classification. Based on the accuracy determine which model should be used.

# Data Collection

---

- Describe how data sets were collected.

Space X API (<https://api.spacexdata.com/v4/rockets/>)

- WebScraping (Wikipedia)

# Data Collection – SpaceX API

---

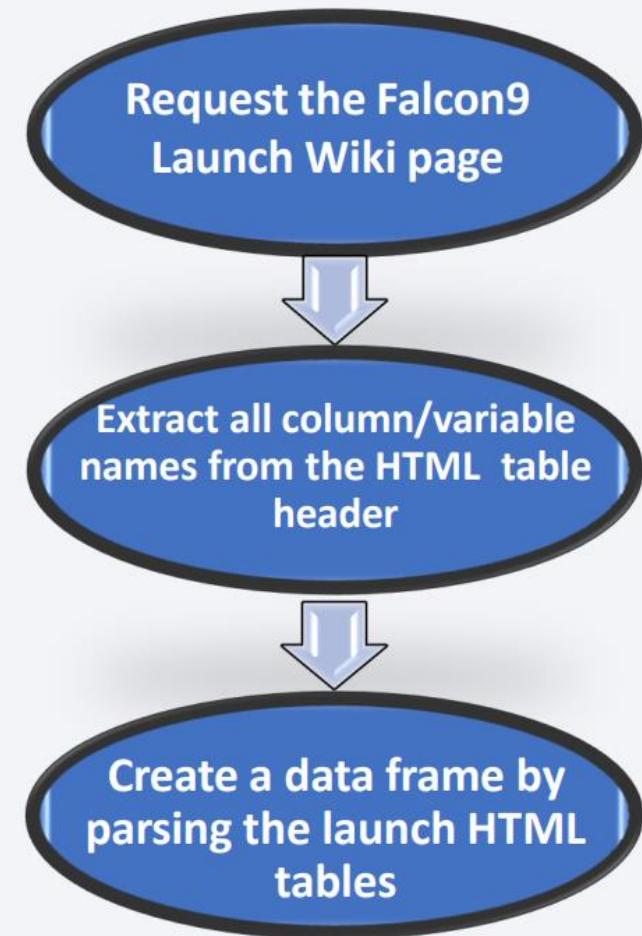
- Present your data collection with SpaceX REST calls using key phrases and flowcharts
- The link to the notebook is  
<https://github.com/SaiKrishnalyer/Data-Structures/blob/main/Code%20of%20DS%20Capstone.pdf>



# Data Collection - Scraping

---

- Data from SpaceX launches can also be obtained from Wikipedia
- **The link to the notebook is [GitHub file for Web scrapping](#)**



# Data Wrangling

---

- Initially some Exploratory Data Analysis (EDA) was performed on the dataset.
- Then the summaries launches per site, occurrences of each orbit and
- occurrences of mission outcome per orbit type were calculated.
- Finally, the landing outcome label was created from Outcome column.

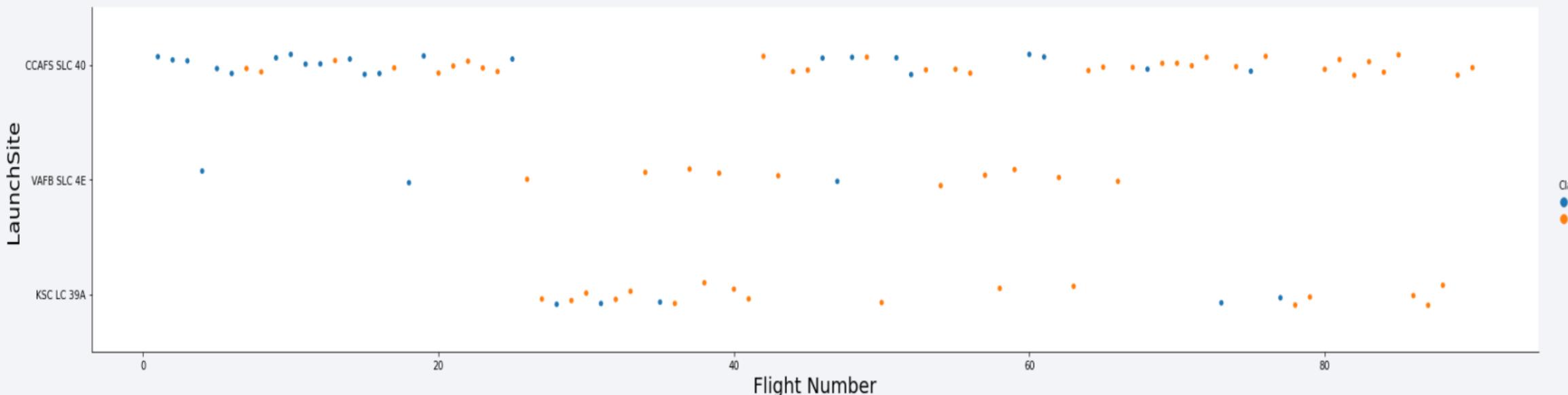


- GitHub File: [Data Wrangling](#)

# EDA with Data Visualization

---

- To explore data, scatterplots and bar plot were used to visualize the relationship between pair of features:
  - Payload Mass X Flight Number, Launch Site X Flight Number, Launch Site X Payload Mass, Orbit and Flight Number, Payload and Orbit



- GitHub File: [EDA Visualization](#)

# EDA with SQL

---

- The following SQL queries were performed:
  - Names of the distinctive launch sites in the space mission;
  - Top 5 launch sites whose name begin with the string 'CCA'
  - Total payload mass carried by boosters launched by NASA (CRS)
  - Average payload mass carried by booster version F9 v1.1
  - Date when the first successful landing outcome in ground pad was achieved
  - Names of the boosters which have success in drone ship and have payload mass between 4000 and 6000 kg
  - Total number of successful and failure mission outcomes
  - Names of the booster versions which have carried the maximum payload mass
  - Failed landing outcomes in drone ship, their booster versions, and launch site names for in year 2015
  - Rank of the count of landing outcomes (such as Failure (drone ship) or Success (ground pad)) between the date 2010-06-04 and 2017-03-20.
- GitHub File: [EDA with SQL](#)

# Build an Interactive Map with Folium

---

- Markers, circles, lines and marker clusters were used with **Folium Maps**
  - Markers indicate points like launch sites
  - Circles indicate highlighted areas around specific coordinates, like NASA Johnson Space Center
  - Marker clusters indicates groups of events in each coordinate, like launches in a launch site
  - Lines are used to indicate distances between two coordinates.
- GitHub File: [Interactive Mapp with Folium](#)

# Build a Dashboard with Plotly Dash

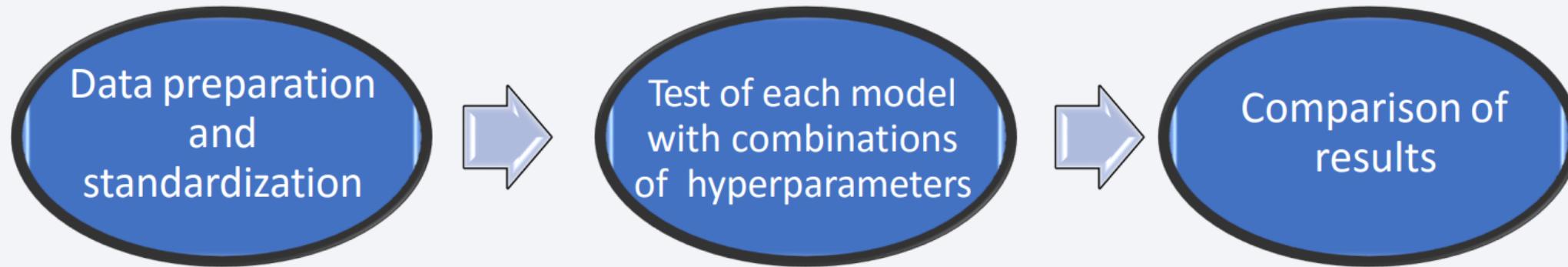
---

- The following graphs and plots were used to visualize data
- Percentage of launches by site
- Payload range
- This combination allowed to quickly analyze the relation between payloads and launch sites, helping to identify where is best place to launch according to payloads.

# Predictive Analysis (Classification)

---

- Four classification models were compared: Logistic regression, SVM(support vector machine), Decision tree and K-nearest neighbors.



- [GitHub File: Predict with ML](#)

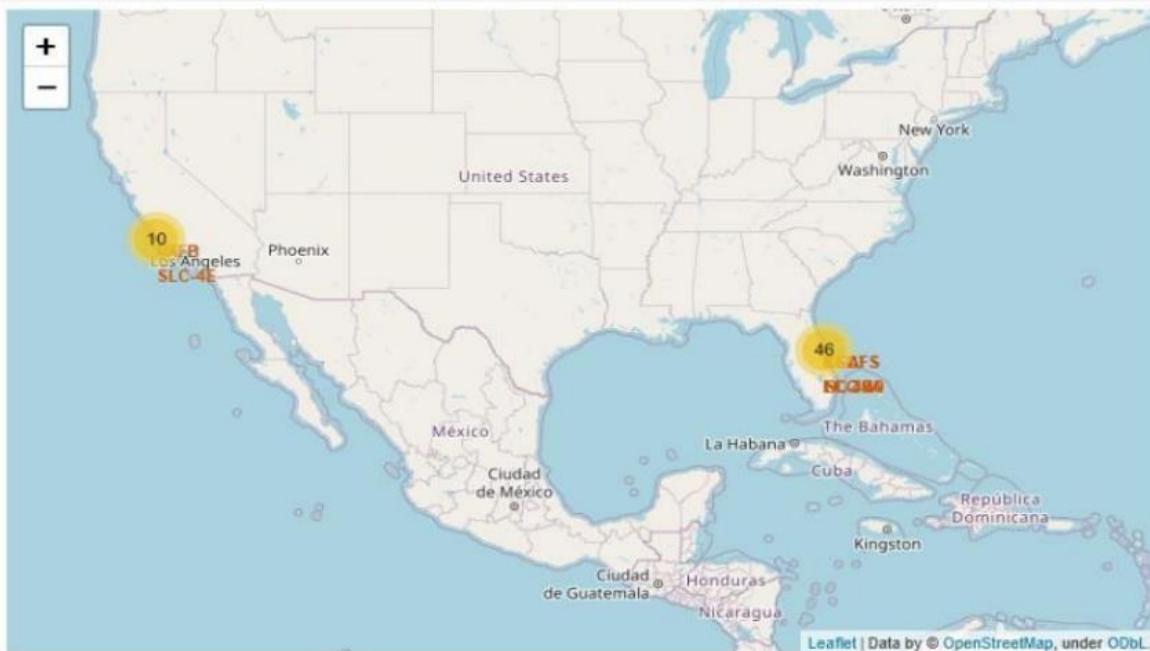
# Results

---

- Exploratory data analysis results:
- Space X uses 4 different launch sites;
- The first launches were done to Space X itself and NASA;
- The average payload of F9 v1.1 booster is 2,928 kg;
- The first success landing outcome happened in 2015 five years after the first launch;
- Many Falcon 9 booster versions were successful at landing in drone ships having payload above the average;
- Almost 100% of mission outcomes were successful;
- Two booster versions failed at landing in drone ships in 2015: F9 v1.1 B1012 and F9 v1.1 B1015;
- The number of landing outcomes became as better as years passed.

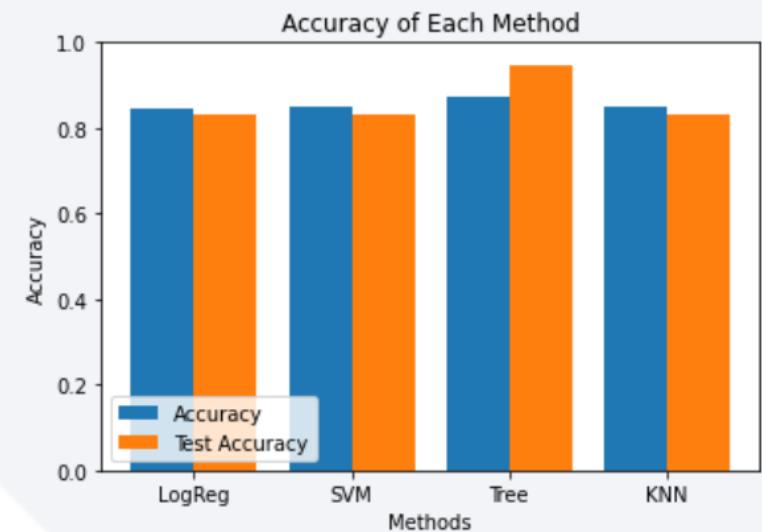
# Results

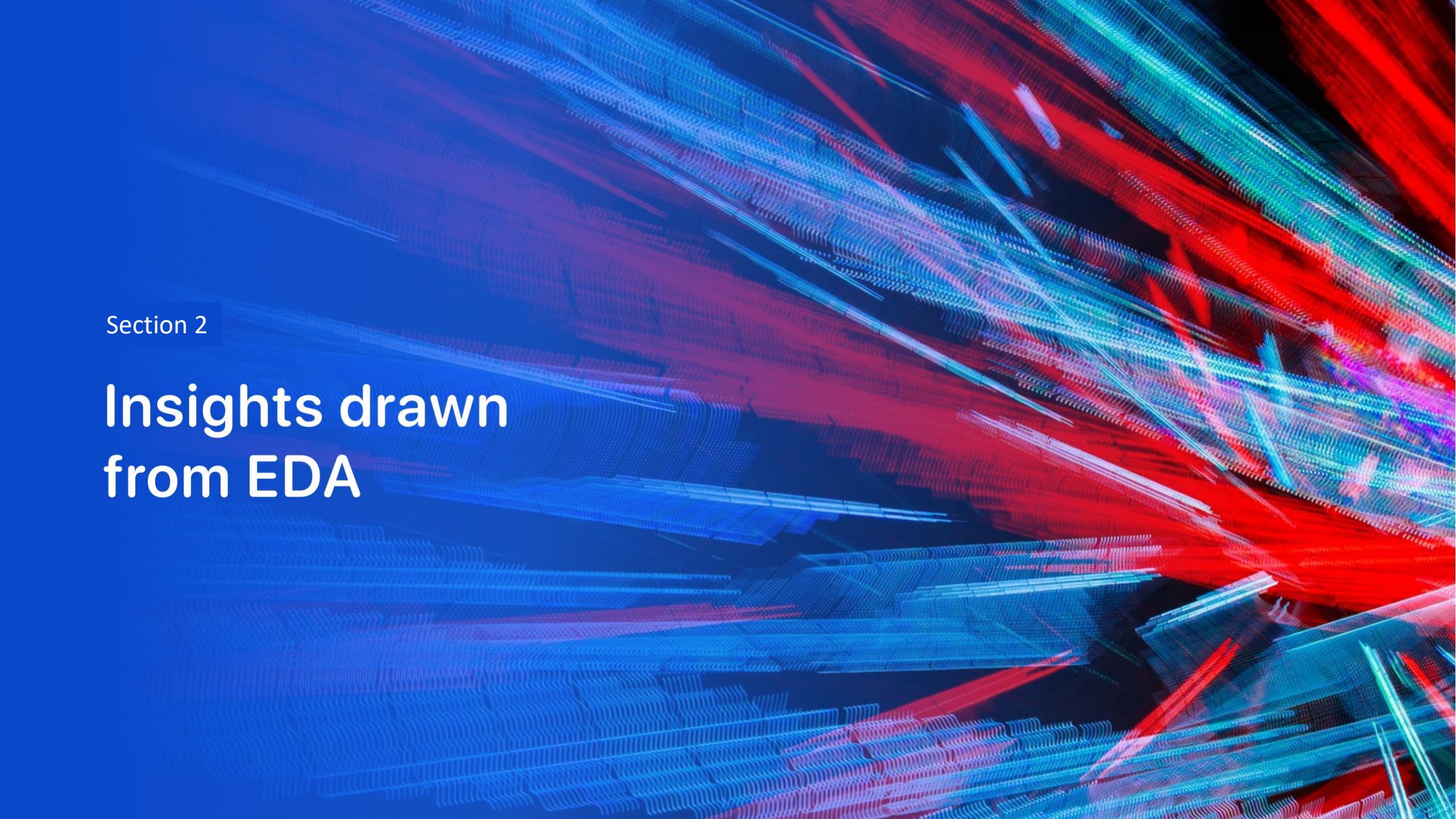
- Using interactive analytics was possible to identify that launch sites use to be in safety places, near sea, for example and have a good logistic infrastructure around.
- Most launches happens at east cost launch sites.



# Results

Predictive Analysis showed that Decision Tree Classifier is the best model to predict successful landings, having accuracy over 87% and accuracy for test data over 94%.

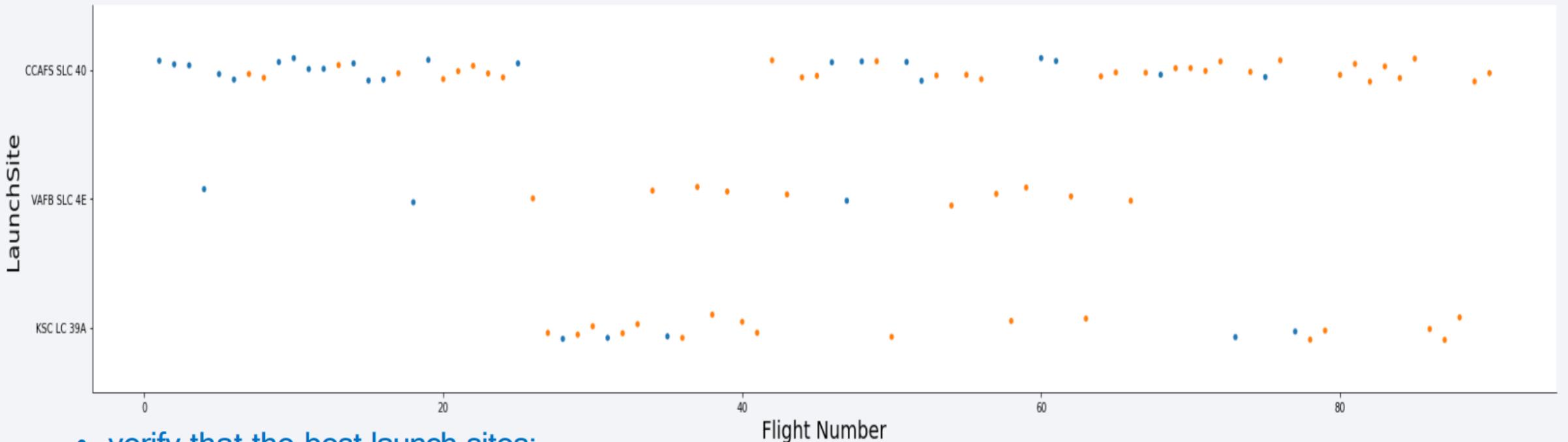


The background of the slide features a complex, abstract digital visualization. It consists of numerous thin, glowing lines that create a sense of depth and motion. The lines are primarily blue and red, with some green and purple highlights. They form a grid-like structure that curves and twists across the frame, resembling a three-dimensional space or a network of data points. The overall effect is futuristic and dynamic.

Section 2

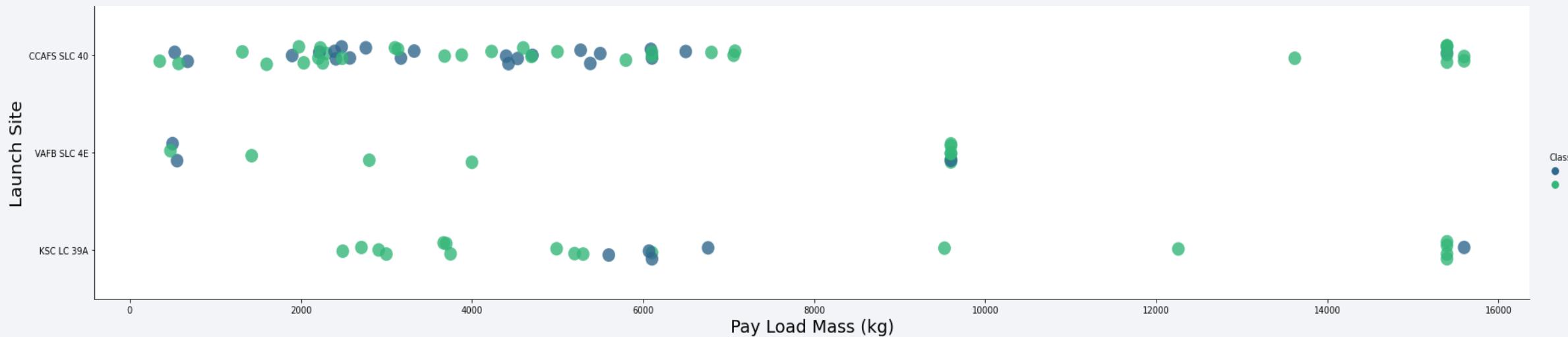
## Insights drawn from EDA

# Flight Number vs. Launch Site



- verify that the best launch sites:
  1. CCAF5 SLC 40
  2. VAFB SLC 4E
  3. KSCLC 39A
- the general success rate improved over time.

# Payload vs. Launch Site

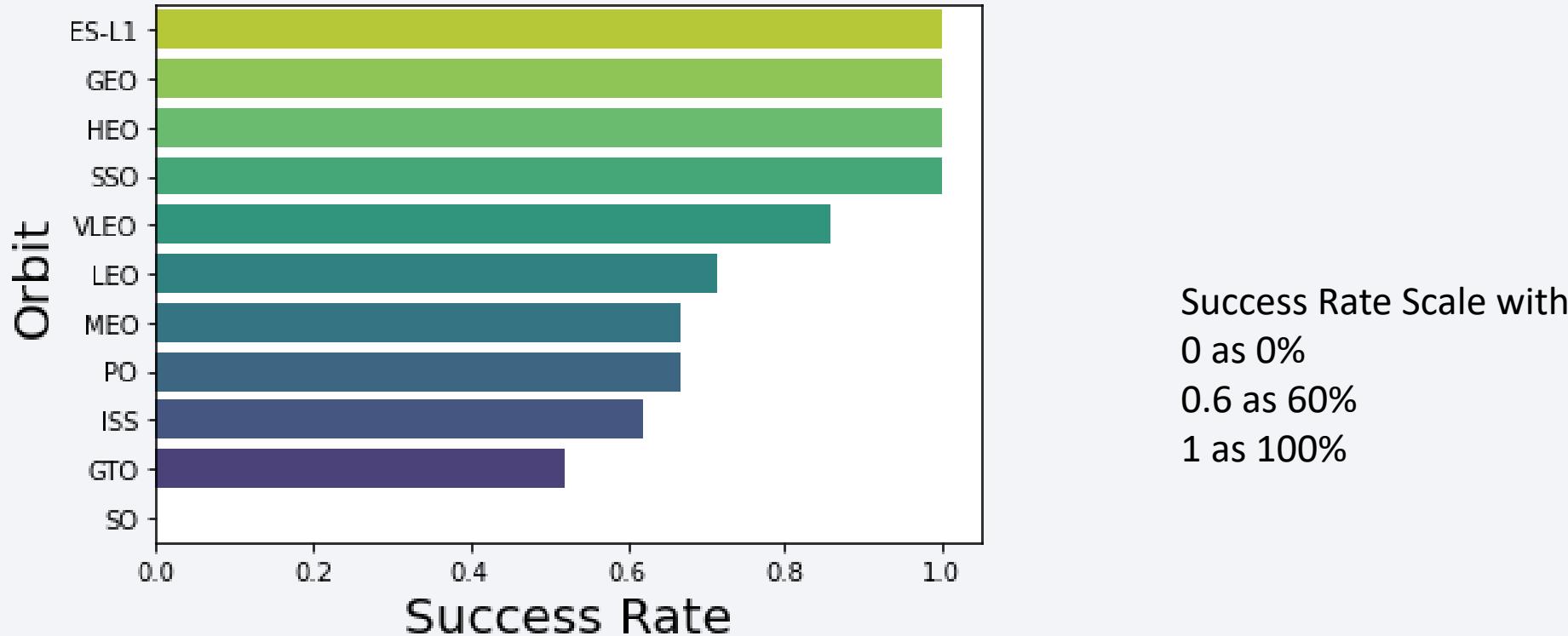


Green indicates successful launch; Purple indicates unsuccessful launch.

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

---



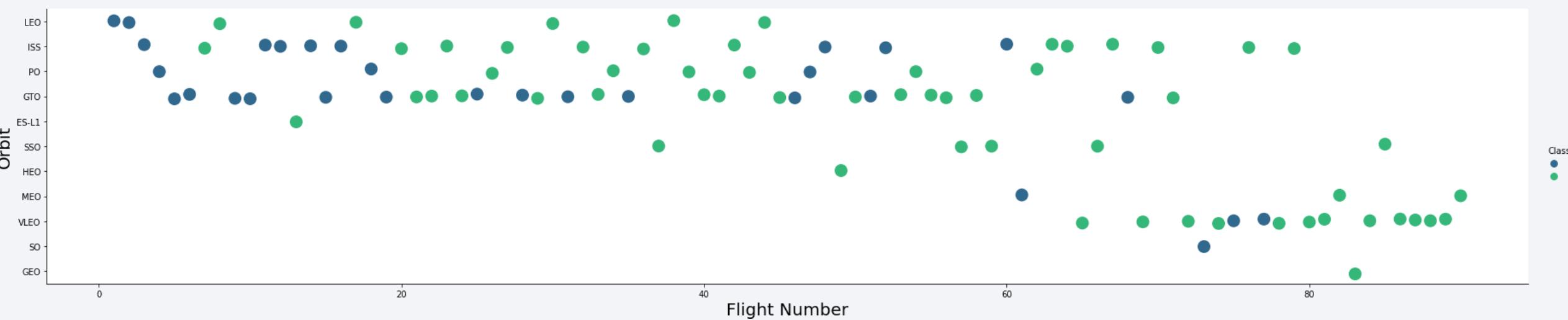
ES-L1 (1), GEO (1), HEO (1) have 100% success rate (sample sizes in parenthesis) SSO (5) has 100% success rate

VLEO (14) has decent success rate and attempts

SO (1) has 0% success rate

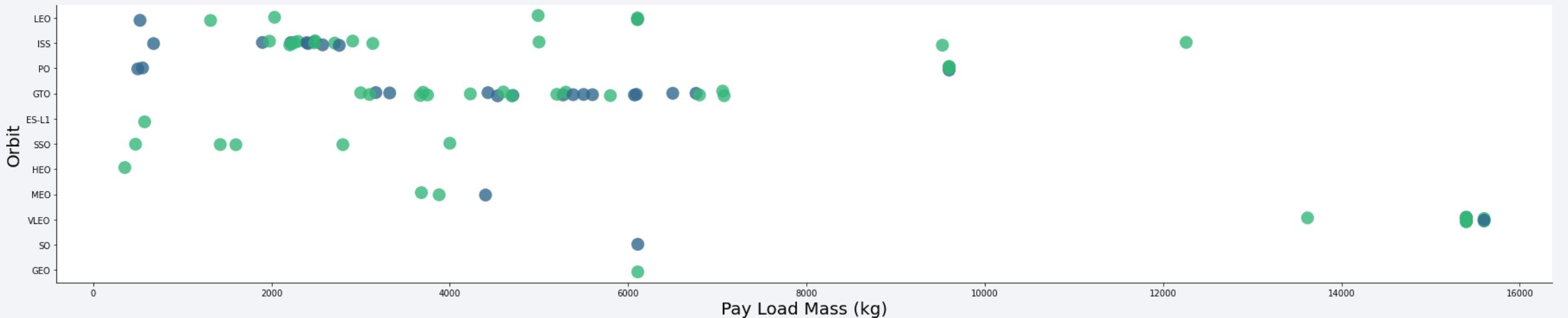
GTO (27) has the around 50% success rate but largest sample

# Flight Number vs. Orbit Type



Green indicates successful launch; Purple indicates unsuccessful launch.

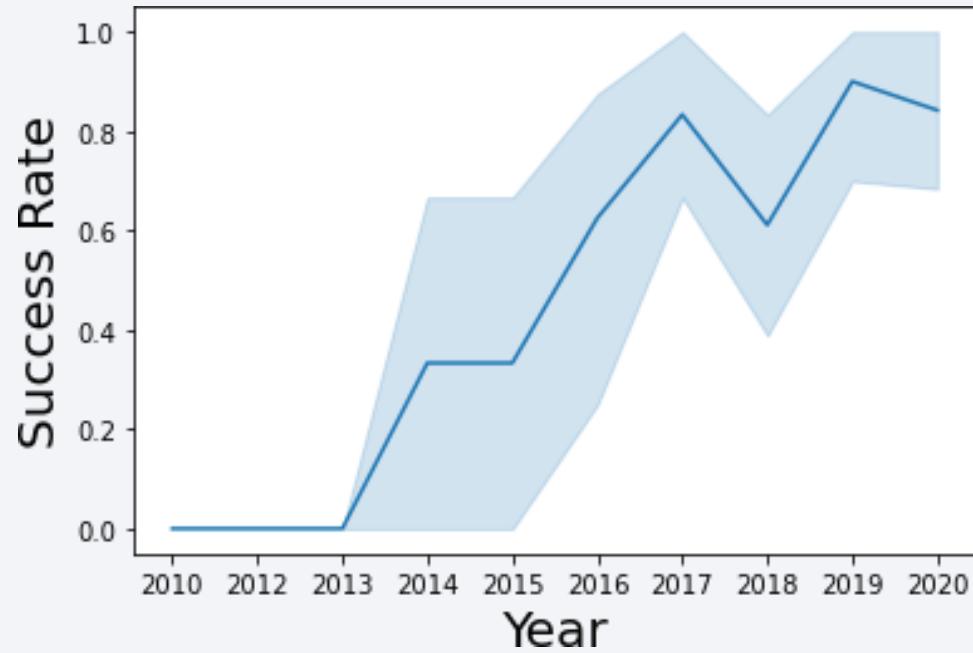
# Payload vs. Orbit Type



Green indicates successful launch; Purple indicates unsuccessful launch.

# Launch Success Yearly Trend

---



95% confidence  
interval (light blue  
shading)

# 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

Query unique launch site names from database.

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'

```
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 beginning with CCA.

# Total Payload Mass

---

```
%%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 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

---

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

<b>first_success</b>
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

---

```
%%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.bs2io90108kqb1od8lcg.firebaseio
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

---

```
%%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

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

---

```
%%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

---

```
%%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.bs2io90108kqb1od8lcg
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.

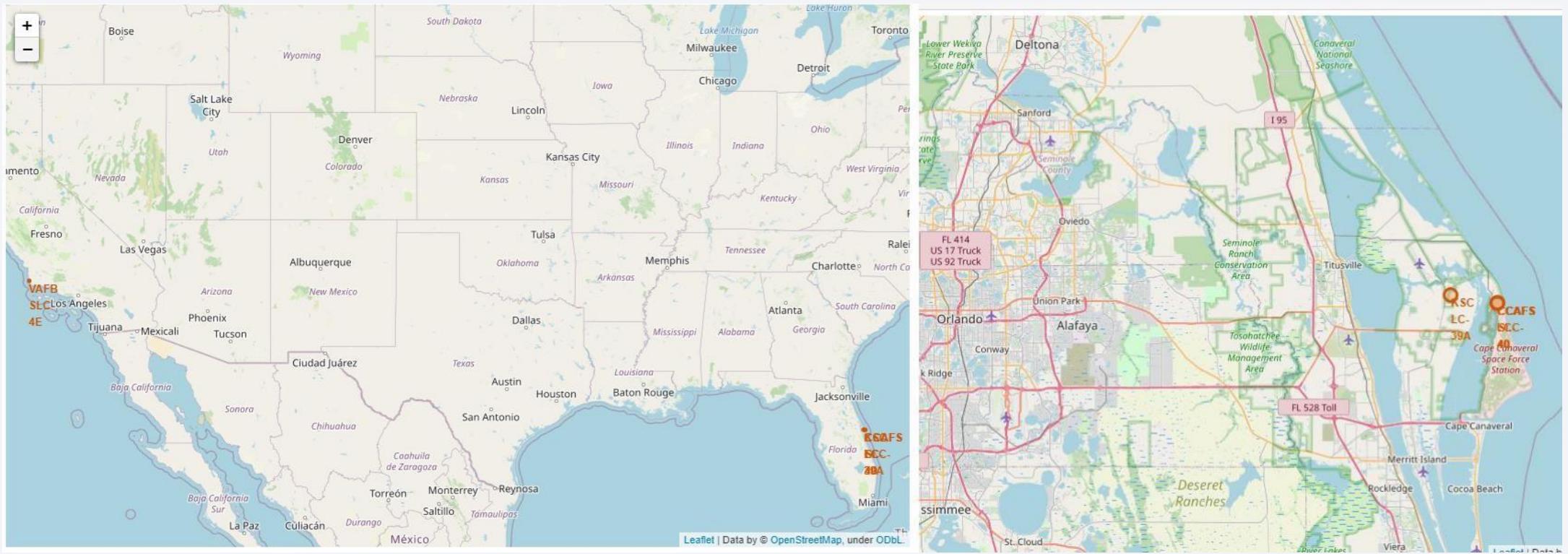
There were 8 successful landings in total during this time period

The background of the slide is a photograph taken from space at night. It shows the curvature of the Earth's horizon against a dark blue sky. Numerous glowing yellow and white points represent city lights, concentrated in coastal and urban areas. In the upper right quadrant, there are bright green and yellow bands of light, likely the Aurora Borealis or Australis. The overall atmosphere is dark and mysterious.

Section 3

# Launch Sites Proximities Analysis

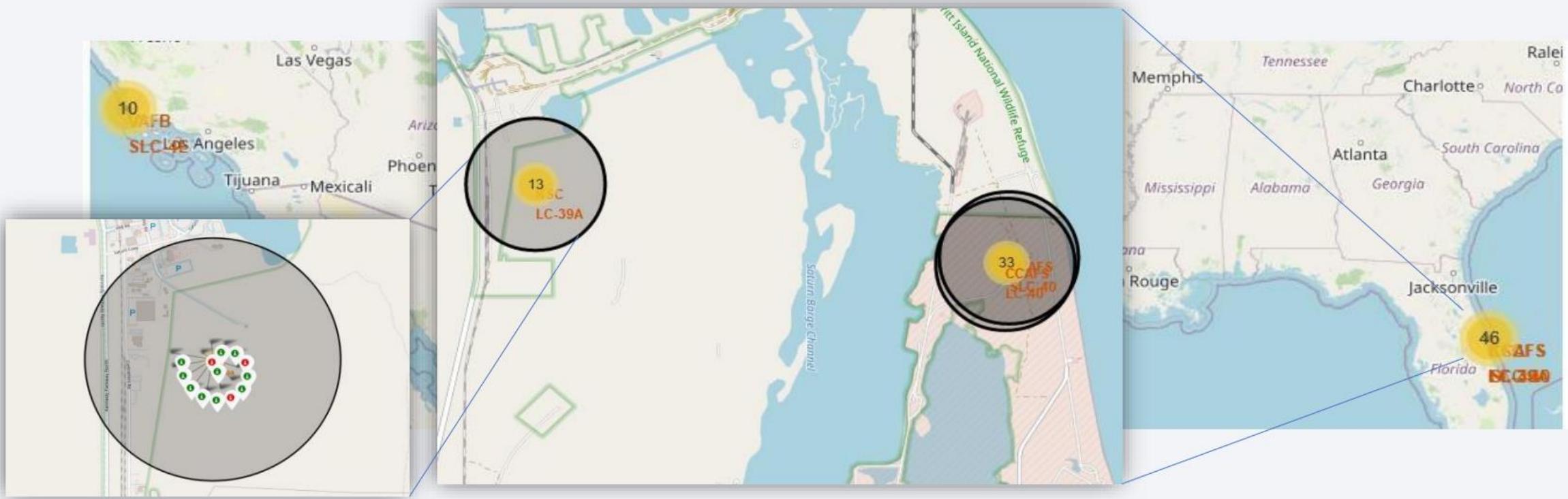
# <Launch Site Locations>



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 Outcomes by Site

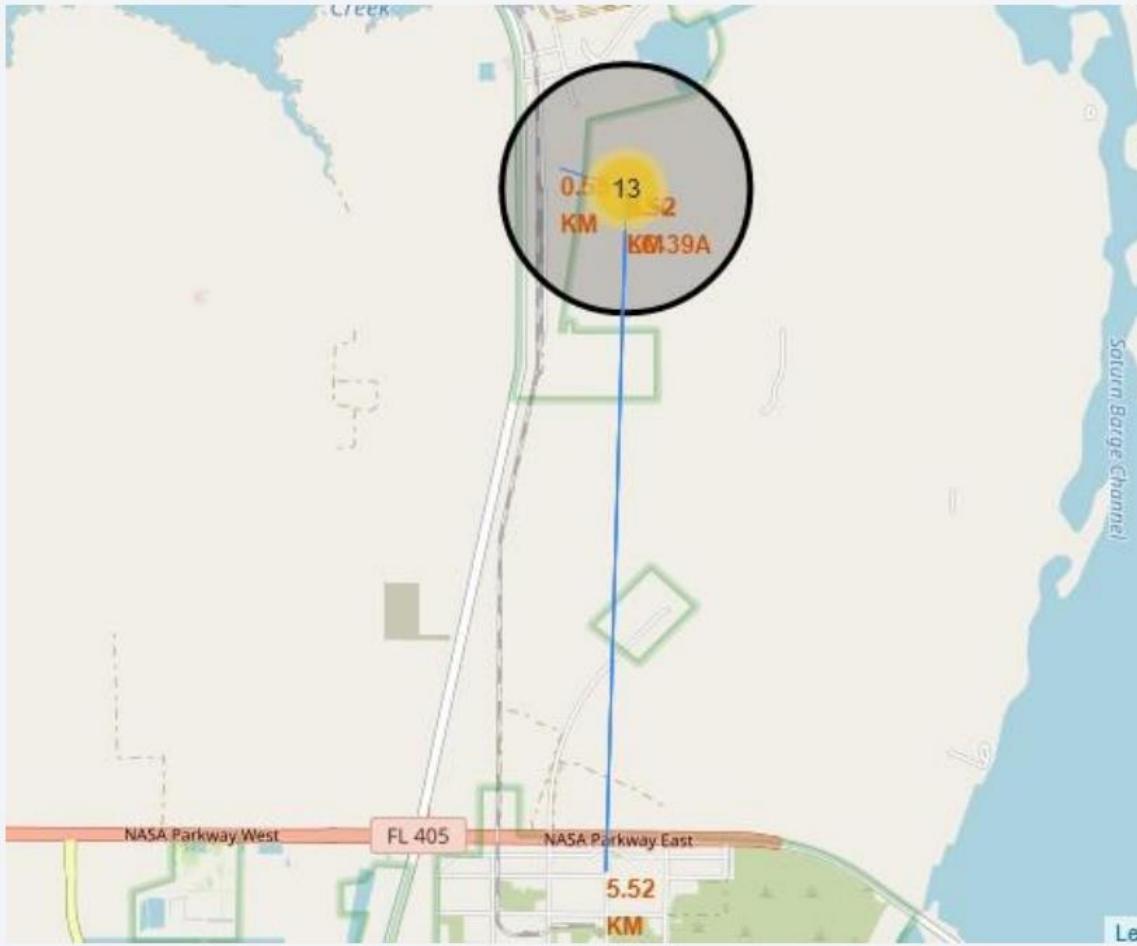
KSC LC-39A launch site launch outcomes



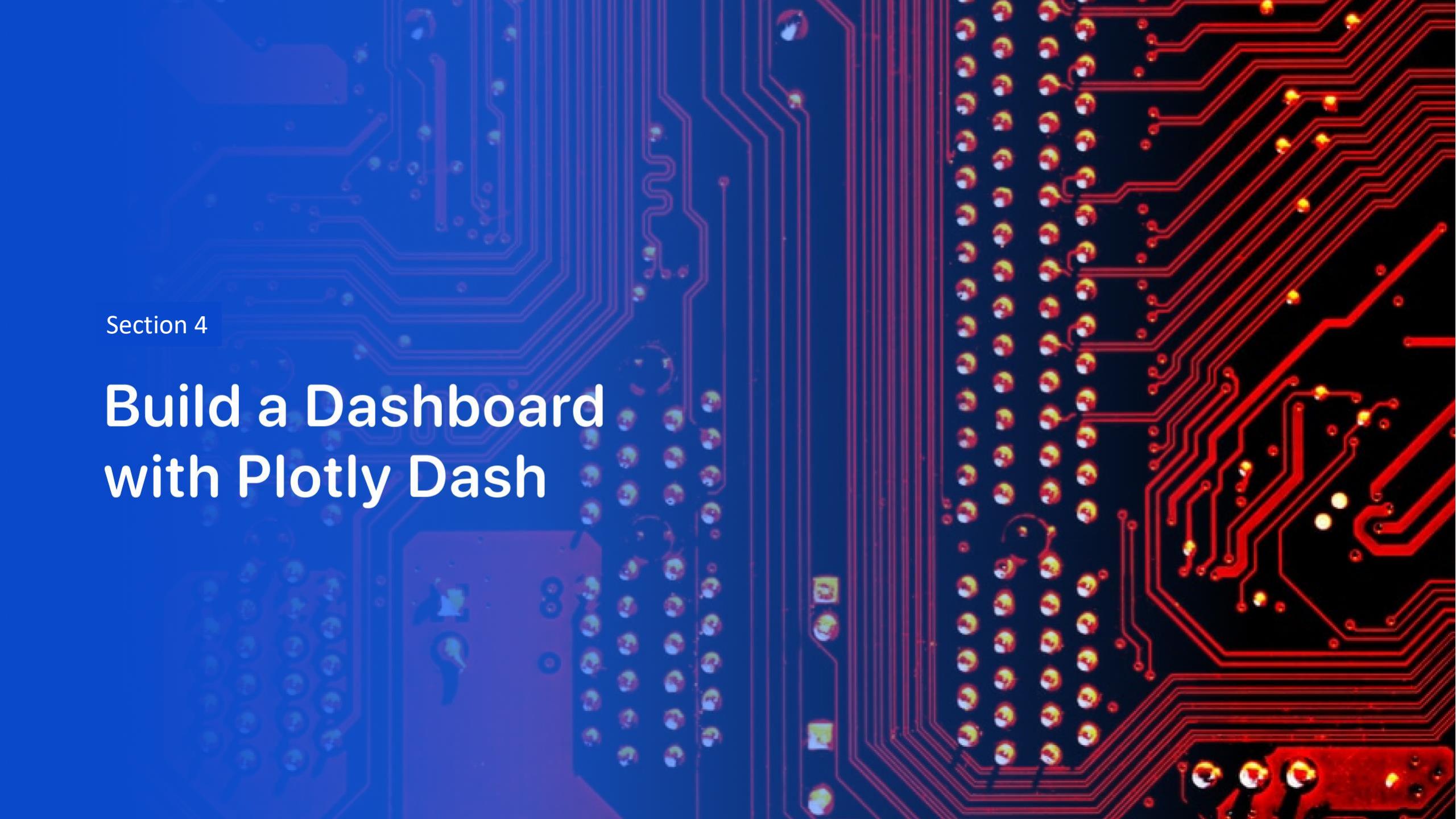
- Green: Successful case
- Red: failure.

# Safety

---



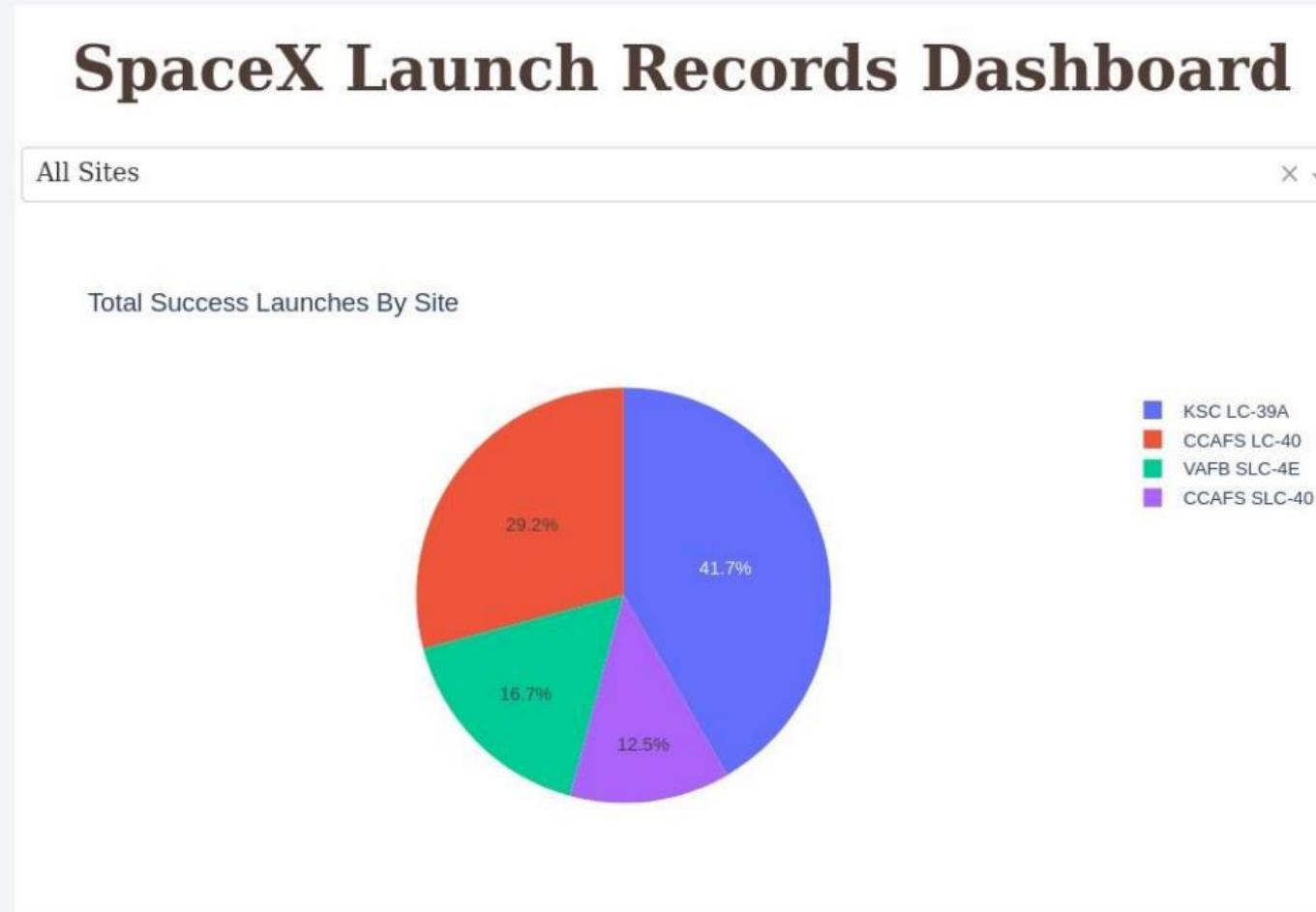
- Site KSCLC-39A has good Safety aspects, near railroad and road and far from inhabited areas (red line).



Section 4

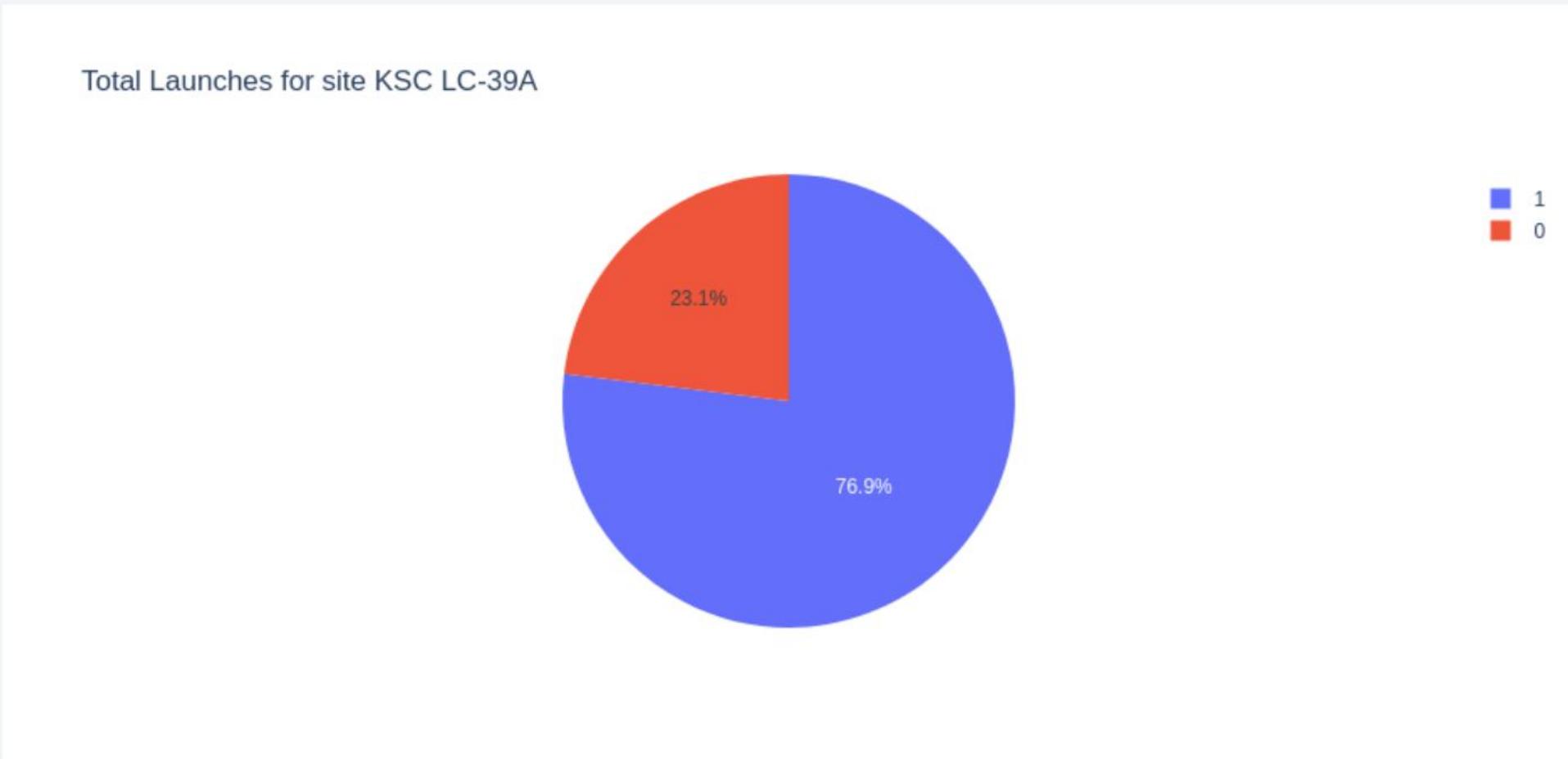
# Build a Dashboard with Plotly Dash

# Successful Launches by Site



# Success Ratio for KSC LC-39A

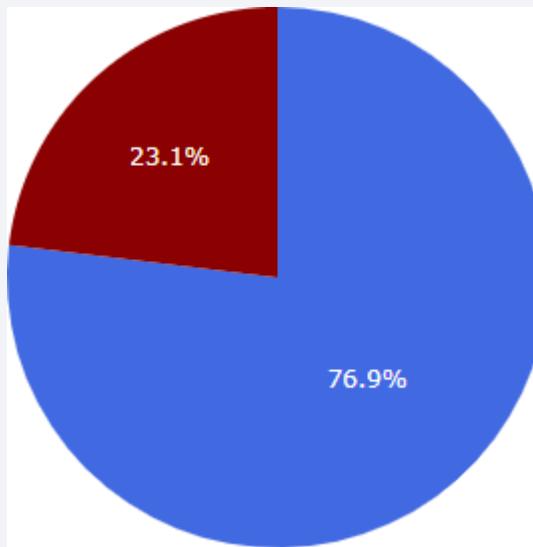
---



- 76.9% of launches are successful.

## Highest Success Rate Launch Site

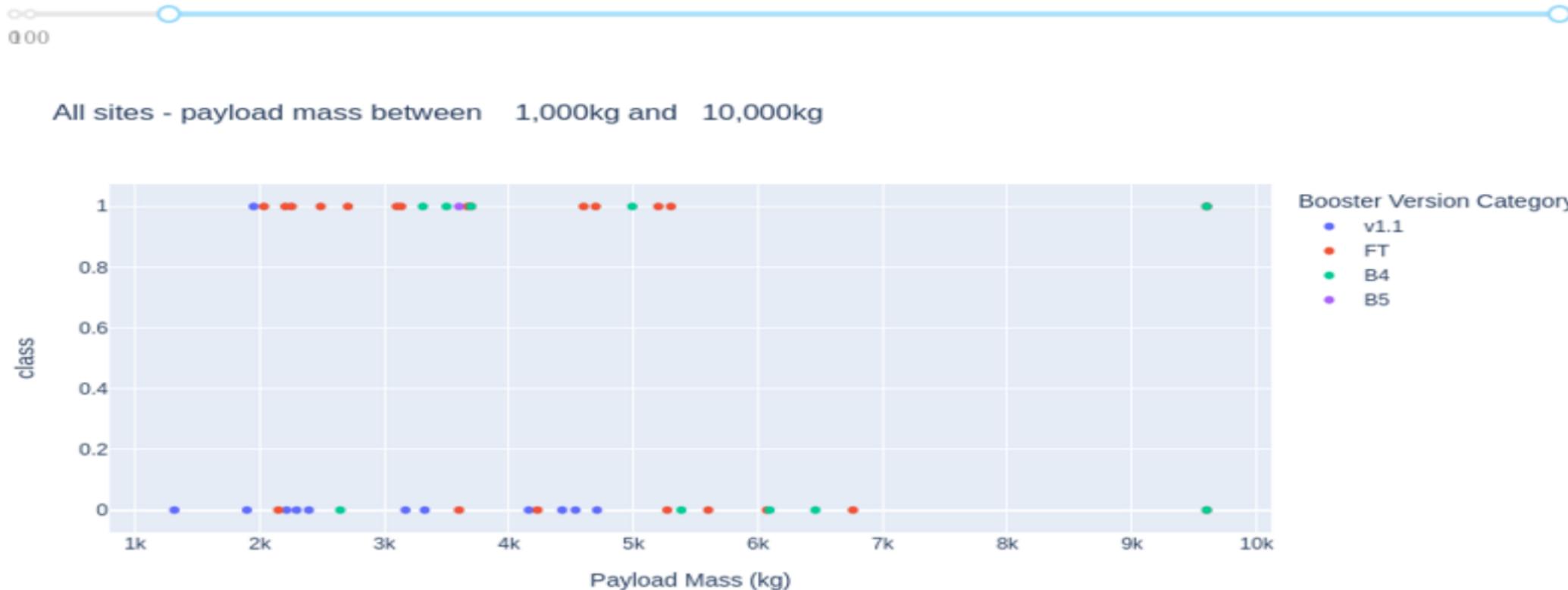
---



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

# Payload vs. Launch Outcome

**Payload range (Kg):**



Payloads under 6,000kg and FT boosters are the most successful combination

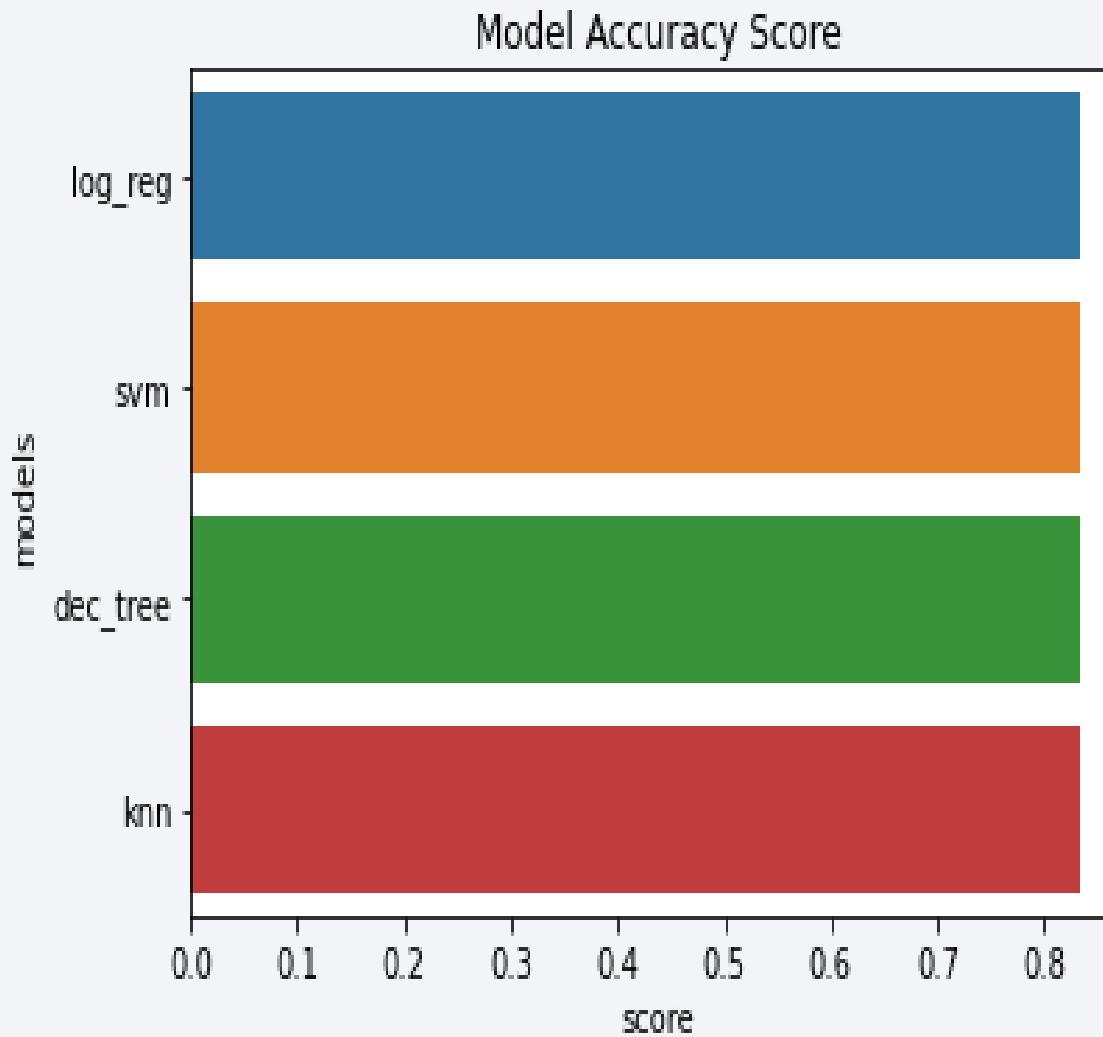
The background of the slide features a dynamic, abstract design. It consists of several thick, curved lines that transition from a bright yellow at the top right to a deep blue at the bottom left. These curves are set against a lighter blue background, creating a sense of motion and depth. The overall effect is reminiscent of a tunnel or a high-speed train track.

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

---

A confusion matrix is a summary of prediction results on a classification problem. The number of correct and incorrect predictions are summarized with count values and broken down by each class. This is the key to the confusion matrix.



# Conclusions

---

The best launch site is K S C L C-39A.

- Launches above 7,000kg have less potential to be failed
- Successful landing outcomes improve over time due to new technologies
- Decision Tree Classifier can be used to predict successful landings and increase profits.

# Appendix

---

- Include any relevant assets like Python code snippets, SQL queries, charts, Notebook outputs, or data sets that you may have created during this project

Instructors: Rav Ahuja, Alex Akison, Aije Egwaikhide, Svetlana Levitan, Romeo Kienzler, Polong Lin, Joseph Santarcangelo, Azim Hirjani, Hima Vasudevan, Saishruthi Swaminathan, Saeed Aghabozorgi, Yan Luo

Special Thanks to All Instructors:

<https://www.coursera.org/professional-certificates/ibm-data-science?#instructors>

- There were Some warning with the codes just ignore them.

Thank you!

