

# Data Science Capstone Project

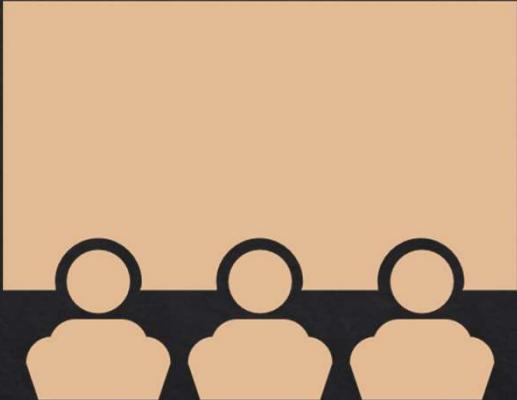
---

Arkadeepto Majumder

1/1/2024

# Outline

---



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

# Executive Summary

---

## Data Collection and Labeling:

- Data was sourced from the public SpaceX API and SpaceX Wikipedia page. The 'class' column was introduced to categorize successful landings. Various techniques, including SQL, visualization, Folium maps, and dashboards, were employed to explore the data. Relevant features were identified for analysis.

## Machine Learning Model Evaluation:

- Four machine learning models—Logistic Regression, Support Vector Machine, Decision Tree Classifier, and K Nearest Neighbors—were created. Despite their similar accuracy rates of approximately 83.33%, all models consistently overpredicted successful landings. The need for additional data is highlighted for more accurate model determination and evaluation.

# Introduction

## Background:

- Commercial Space Age is Here
- Space X has best pricing (\$62 million vs. \$165 million USD)
- Largely due to ability to recover part of rocket (Stage 1)
- Space Y wants to compete with Space X

## Problem:

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



SpaceX Falcon 9 Rocket – The Verge

# Methodology

OVERVIEW OF DATACOLLECTION, WRANGLING, VISUALIZATION,  
DASHBOARD, AND MODEL METHODS

# Methodology

---

- Data collection methodology:
  - Combined data from SpaceX public API and SpaceX Wikipedia page
- Perform data wrangling
  - 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 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.

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.

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

GitHub url:

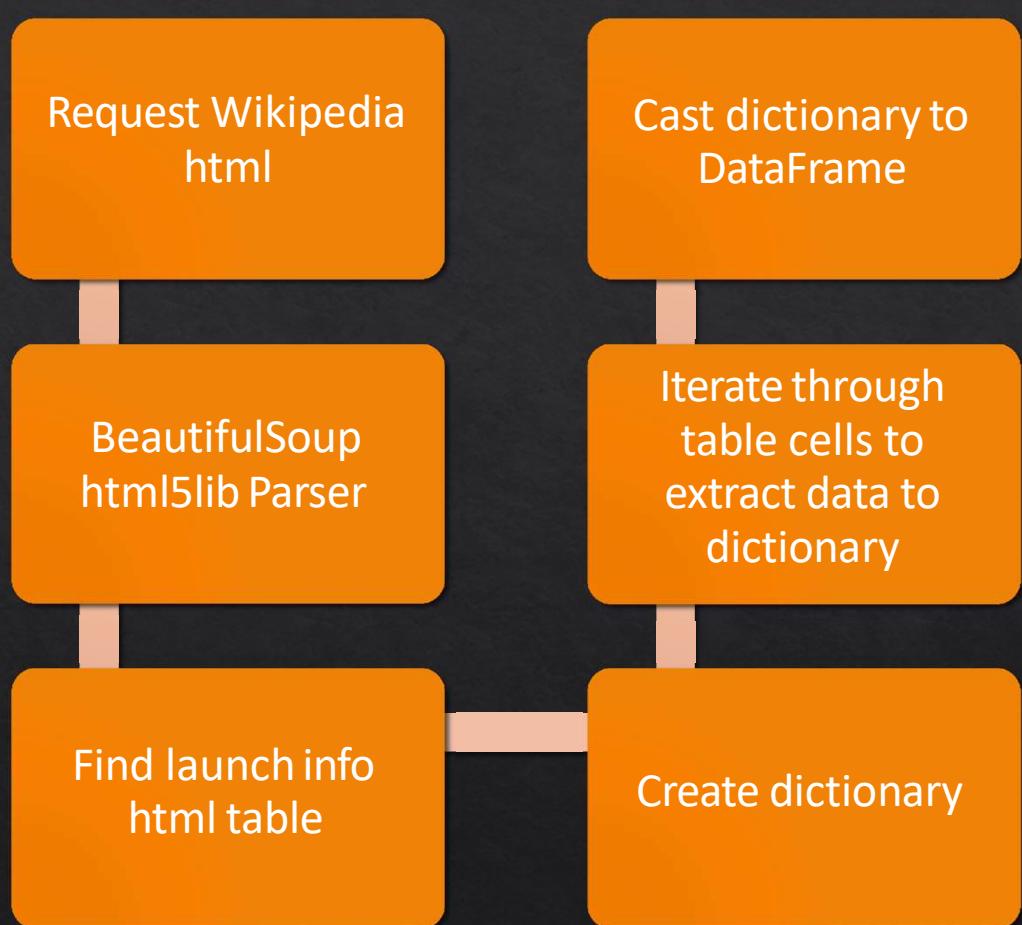
<https://github.com/Arkadeepto2002/IBM-Data-Science-Professional-Certificate/blob/main/Applied%20Data%20Science%20Capstone/Data%20Collection%20Api.ipynb>



# Data Collection – Web Scraping

GitHub url:

<https://github.com/Arkadeepto2002/IBM-Data-Science-Professional-Certificate/blob/main/Applied%20Data%20Science%20Capstone/Data%20Collection%20with%20Web%20Scraping.ipynb>



# Data Wrangling

- 
- ❖ Create 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/Arkadeepto2002/IBM-Data-Science-Professional-Certificate/blob/main/Applied\\_Data\\_Science\\_Capstone/Data%20wrangling.ipynb](https://github.com/Arkadeepto2002/IBM-Data-Science-Professional-Certificate/blob/main/Applied_Data_Science_Capstone/Data%20wrangling.ipynb)

# EDA with Visualization

EXPLORATORY DATA ANALYSIS WITH SEABORN PLOTS

# EDA with Data Visualization

---

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

Plots Used:

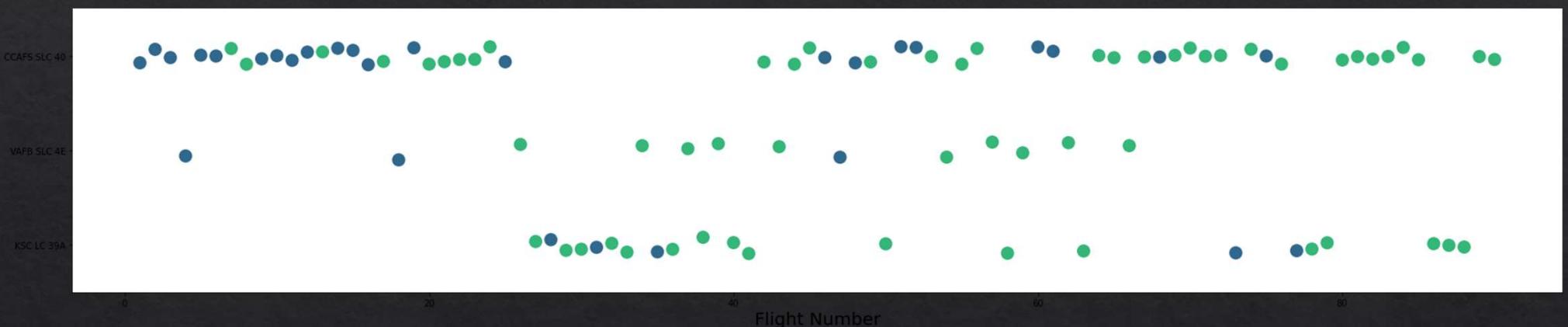
Flight Number vs. Payload Mass, Flight Number vs. Launch Site, Payload Mass vs. Launch Site, Orbit vs. Success Rate, Flight Number vs. Orbit, Payload vs Orbit, and Success Yearly Trend

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

GitHub url:

[https://github.com/Arkadeepto2002/IBM-Data-Science-Professional-Certificate/blob/main/Applied\\_Data\\_Science\\_Capstone/EDA%20with%20Visualization.ipynb](https://github.com/Arkadeepto2002/IBM-Data-Science-Professional-Certificate/blob/main/Applied_Data_Science_Capstone/EDA%20with%20Visualization.ipynb)

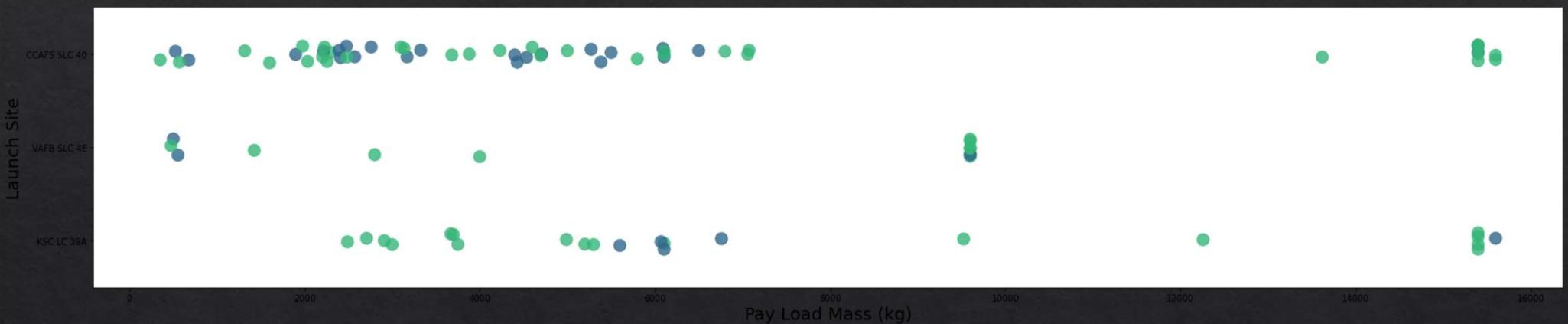
# Flight Number vs. LaunchSite



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

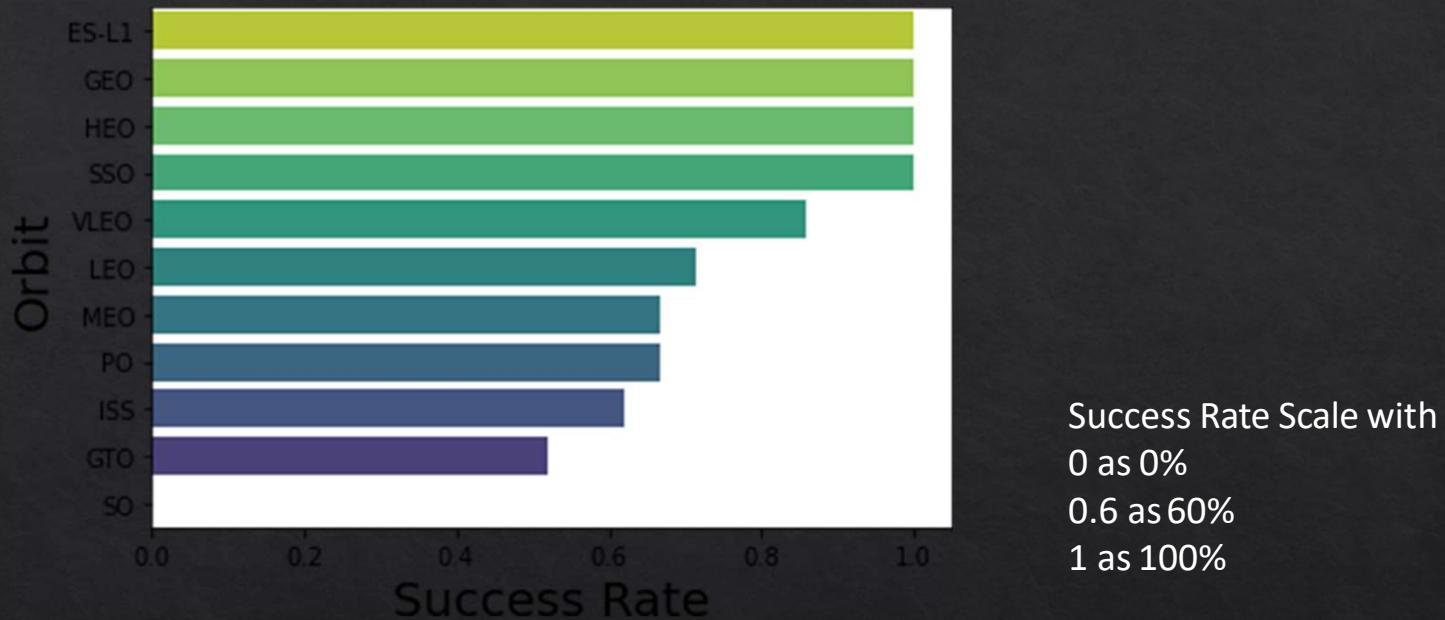


Green indicates successful launch; Purple indicates unsuccessful launch.

Payload mass appears to fall mostly between 0-6000 kg.

Different launch sites also seem to use different payload mass.

## Successrate vs. Orbittype



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

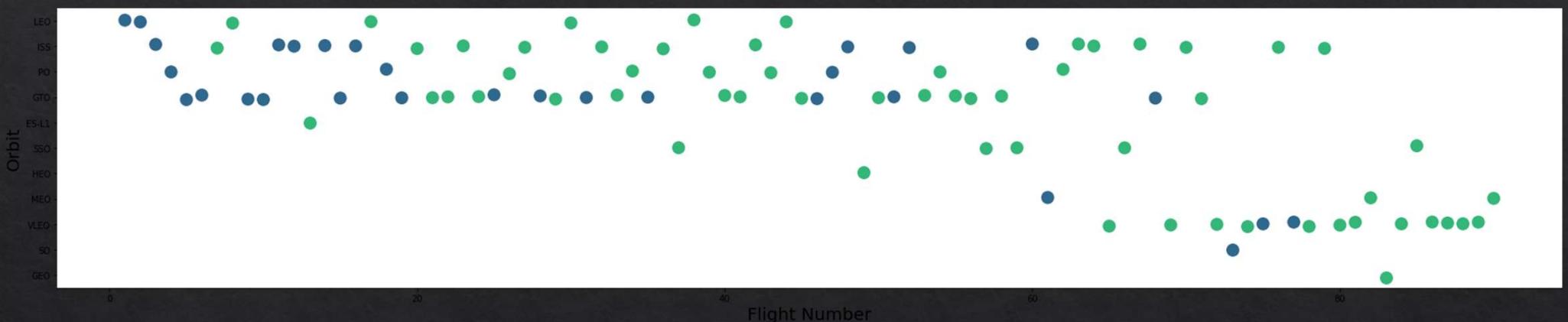
SSO (5) has 100% success rate

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. Orbittype



Green indicates successful launch; Purple indicates unsuccessful launch.

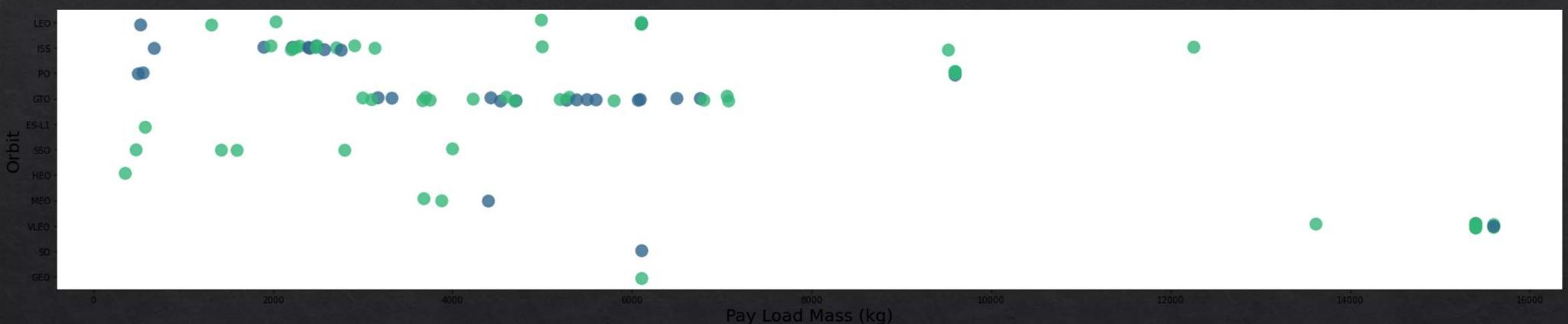
Launch Orbit preferences changed over Flight Number.

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. Orbittype



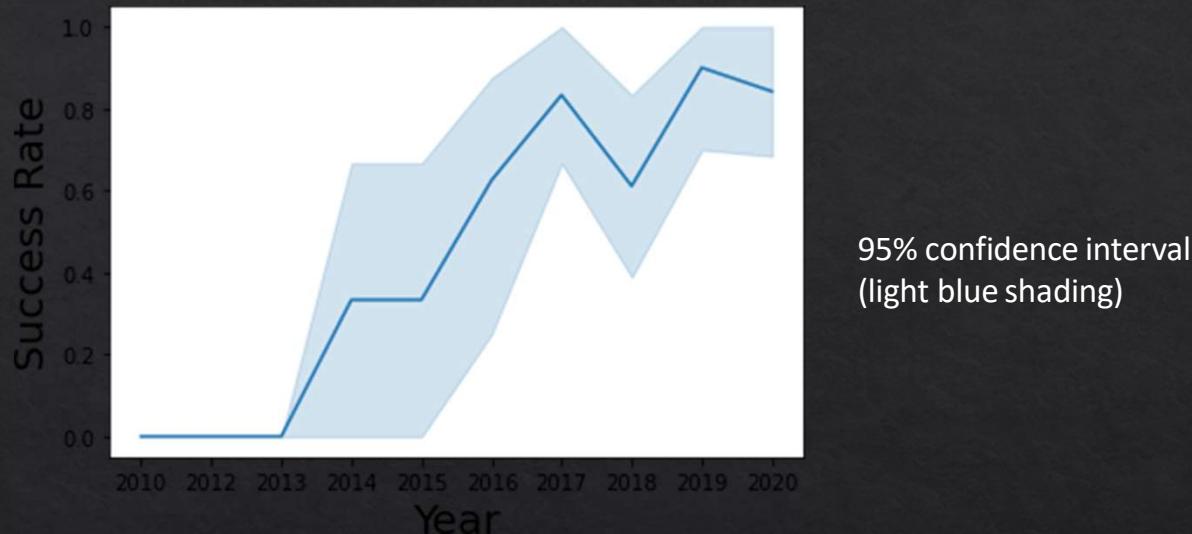
Green indicates successful launch; Purple indicates unsuccessful launch.

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%

# EDAwithSQL

EXPLORATORY DATA ANALYSIS WITH SQL DB2  
INTEGRATED IN PYTHON WITH SQLALCHEMY

# 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/Arkadeepto2002/IBM-Data-Science-Professional-Certificate/blob/main/Applied\\_Data\\_Science\\_Capstone/EDA%20with%20SQL.ipynb](https://github.com/Arkadeepto2002/IBM-Data-Science-Professional-Certificate/blob/main/Applied_Data_Science_Capstone/EDA%20with%20SQL.ipynb)

# All Launch Site Names

```
In [4]: %%sql  
SELECT UNIQUE LAUNCH_SITE  
FROM SPACEXDATASET;  
  
* ibm_db_sa://ftb12020:***@0c77d6f2  
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 Beginning with 'CCA'

| In [5]:    | %%sql<br>SELECT *<br>FROM SPACEXDATASET<br>WHERE LAUNCH_SITE LIKE 'CCA%'<br>LIMIT 5;  |                 |             |   |                 |           |                    |                 |                     |
|------------|---|-----------------|-------------|---|-----------------|-----------|--------------------|-----------------|---------------------|
| Out[5]:    | * ibm_db_sa://ftb12020:***@0c77d6f2-5da9-48a9-81f8-86b520b87518.bs2io90108kqb1od8lcg.databases.appdomain.cloud:31198/bludb<br>Done. |                 |             |   |                 |           |                    |                 |                     |
| 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)<br>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 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 F9v1.1

```
%%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 Pad 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    |

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.bs2io90108kqb1od81cg.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 Each Mission Outcome

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

| 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 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.bs2io90l08kqb1od81cg.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.

# 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 'Succes%' 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

# Interactive Map with Folium

# Building an interactive map with Folium

---

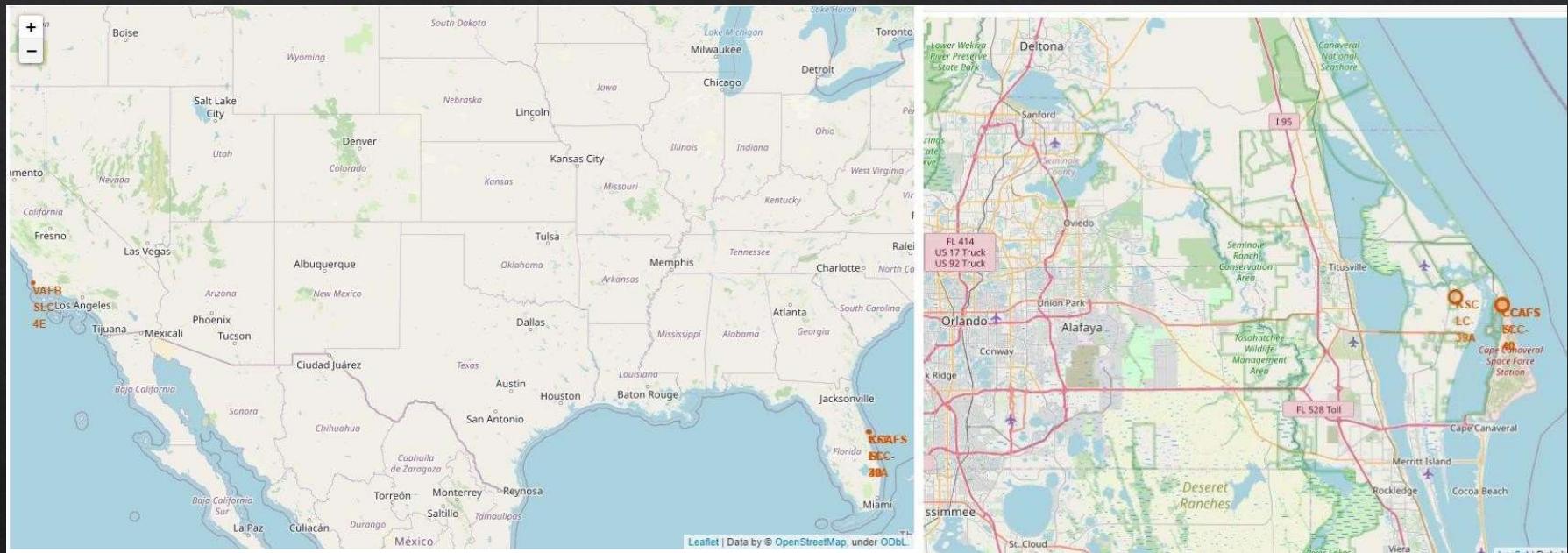
Folium maps mark Launch Sites, successful and unsuccessful landings, and a proximity example to key locations: Railway, Highway, Coast, and City.

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

[GitHub url:](#)

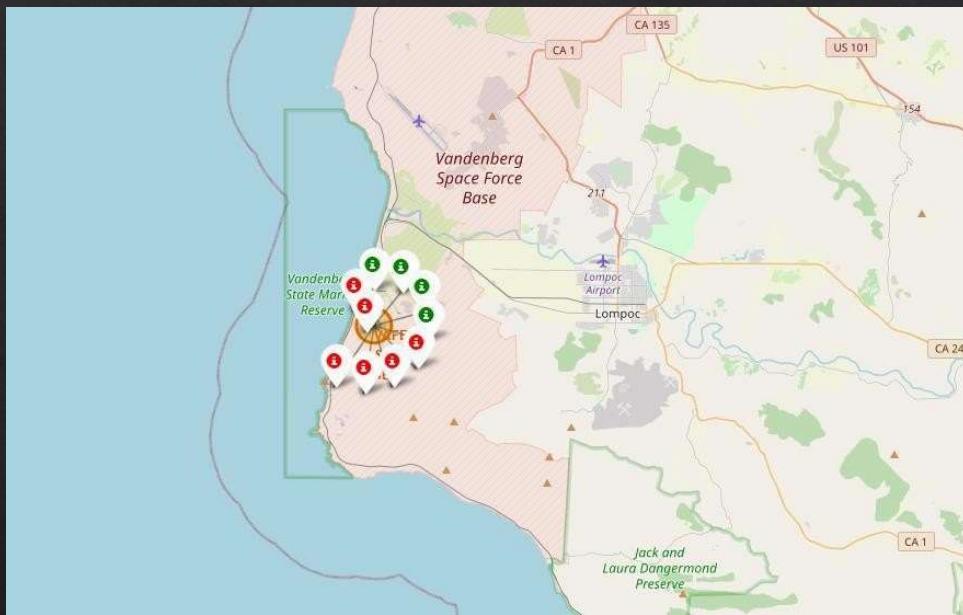
[https://github.com/Arkadeepto2002/IBM-Data-Science-Professional-Certificate/blob/main/Applied\\_Data\\_Science\\_Capstone/Interactive%20Dashboard.ipynb](https://github.com/Arkadeepto2002/IBM-Data-Science-Professional-Certificate/blob/main/Applied_Data_Science_Capstone/Interactive%20Dashboard.ipynb)

# Launch Site Locations

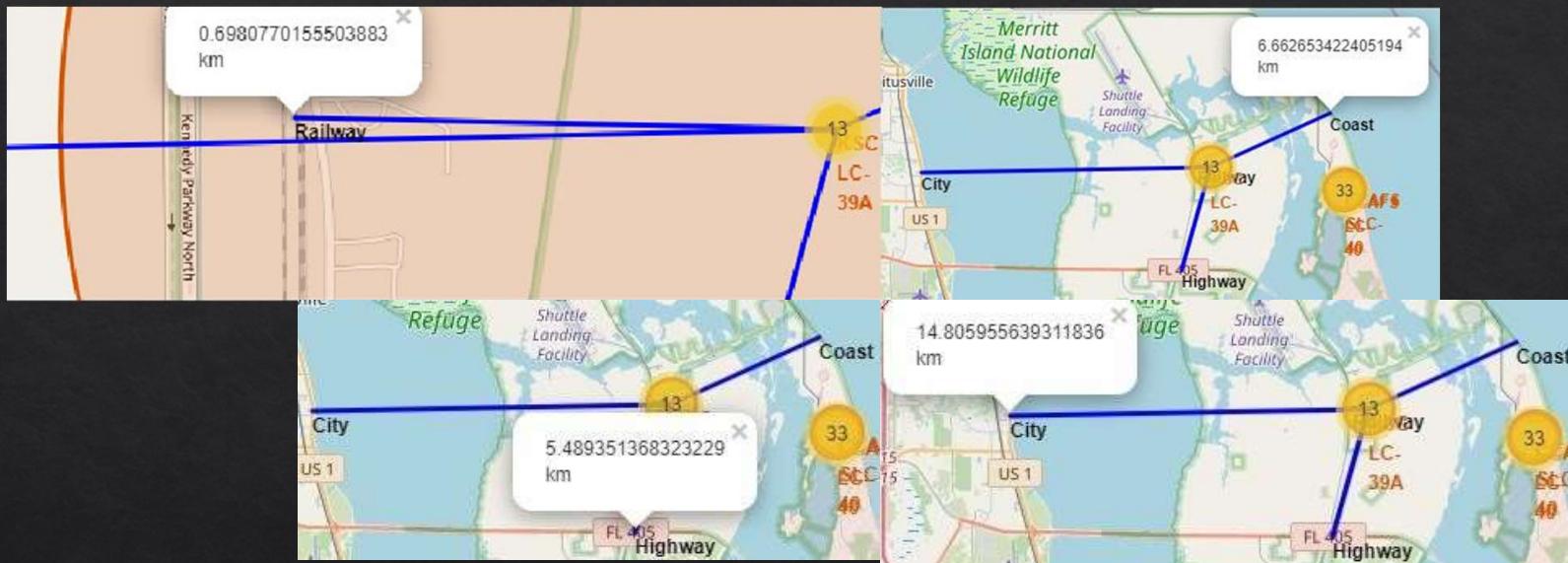


# Color-Coded Launch Markers

---



# Key Location Proximities



# Build a Dashboard with PlotlyDash

---

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 10000 kg.

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

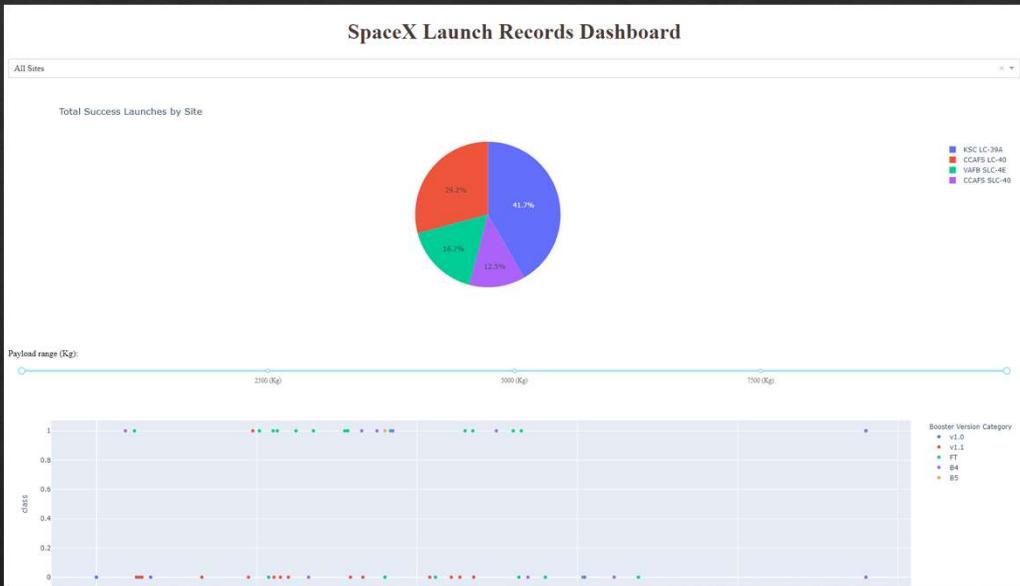
The scatter plot can help us see how success varies across launch sites, payload mass, and booster version category.

GitHub url:

[https://github.com/Arkadeepto2002/IBM-Data-Science-Professional-Certificate/blob/main/Applied\\_Data\\_Science\\_Capstone/spacex\\_dash\\_app.py](https://github.com/Arkadeepto2002/IBM-Data-Science-Professional-Certificate/blob/main/Applied_Data_Science_Capstone/spacex_dash_app.py)

# Build a Dashboard with Plotly Dash

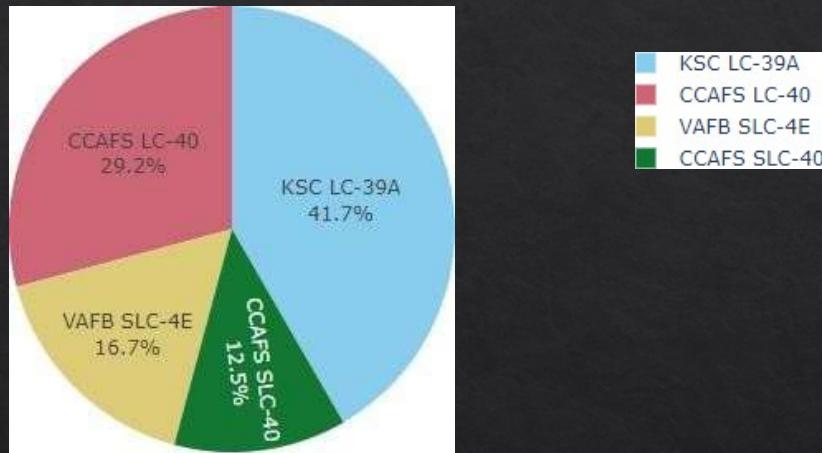
# Results



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.

# Successful Launches Across Launch Sites

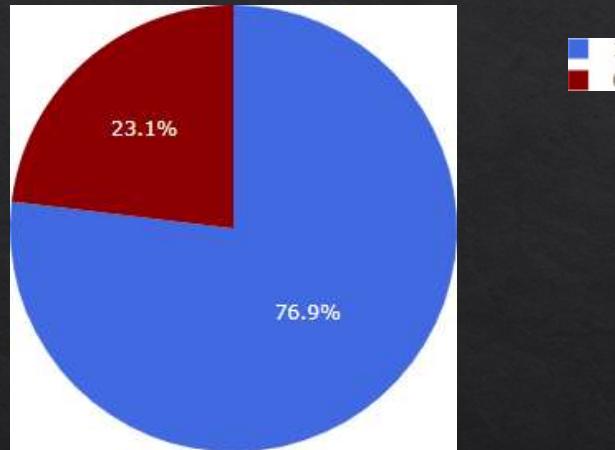
---



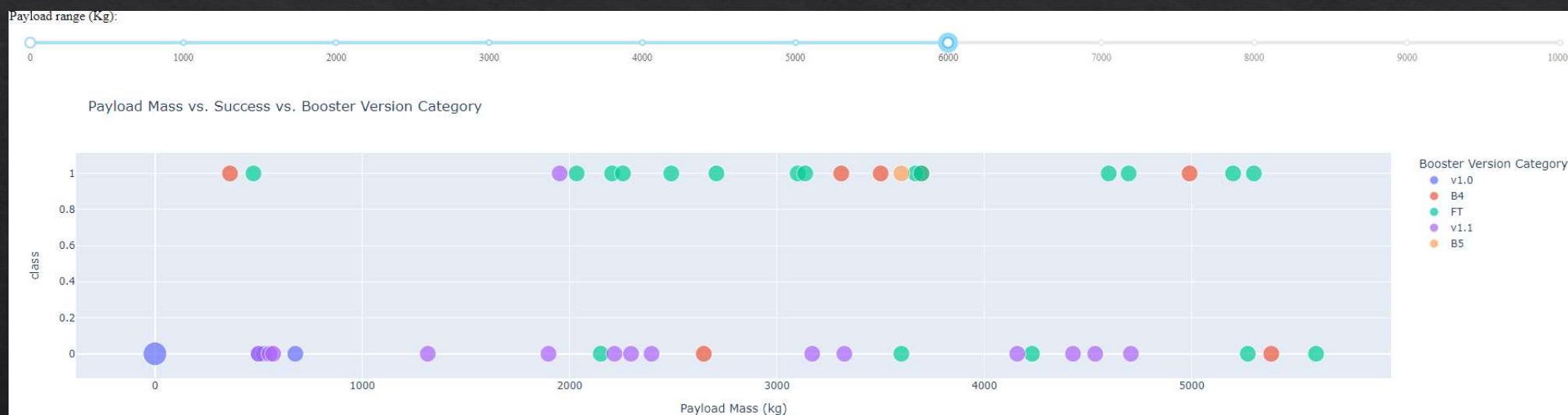
# Highest Success Rate Launch Site

---

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



# Payload Mass vs. Success vs. Booster Version Category



# Predictive Analysis

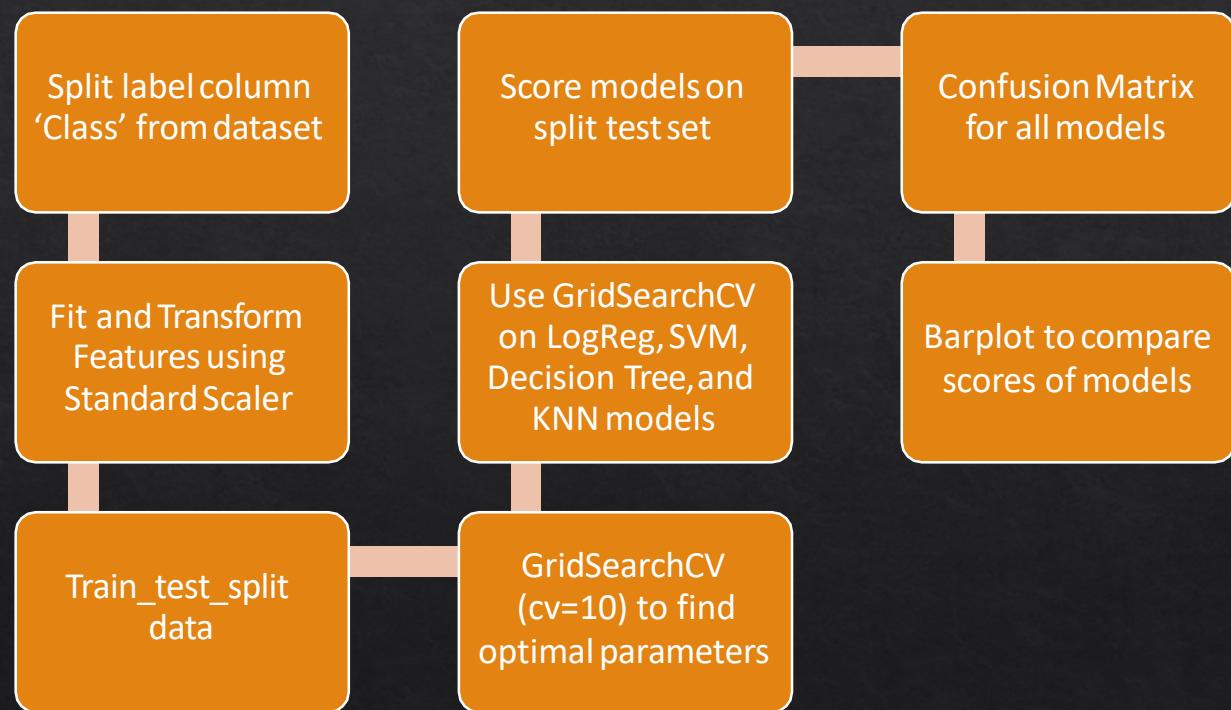
---

GRIDSEARCHCV(CV=10) ON LOGISTIC REGRESSION, SVM, DECISION TREE, AND KNN

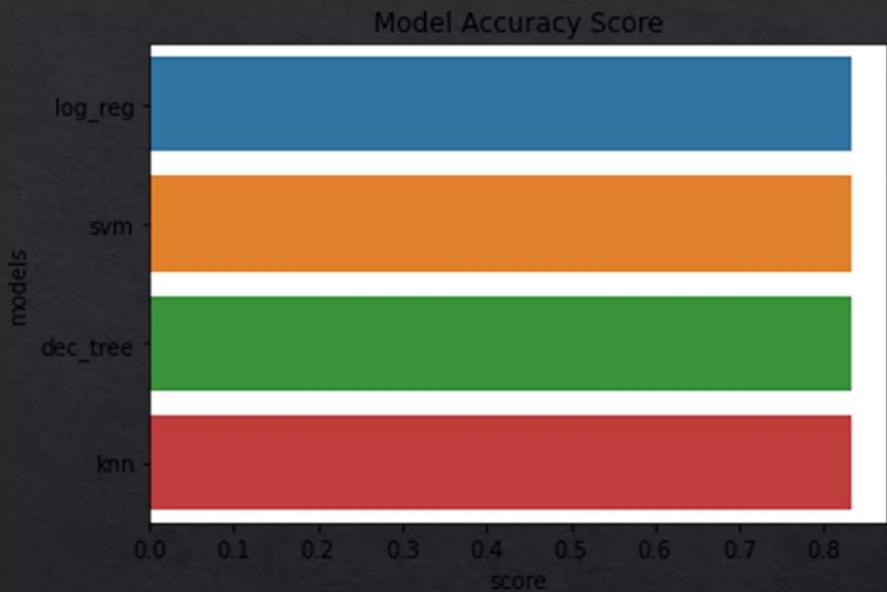
# Predictive analysis(Classification)

GitHub url:

[https://github.com/Arkadeep-to2002/IBM-Data-Science-Professional-Certificate/blob/main/Applying\\_Data\\_Science\\_Capstone/Machine%20learning%20models%20to%20predict%20the%20outcome.ipynb](https://github.com/Arkadeep-to2002/IBM-Data-Science-Professional-Certificate/blob/main/Applying_Data_Science_Capstone/Machine%20learning%20models%20to%20predict%20the%20outcome.ipynb)

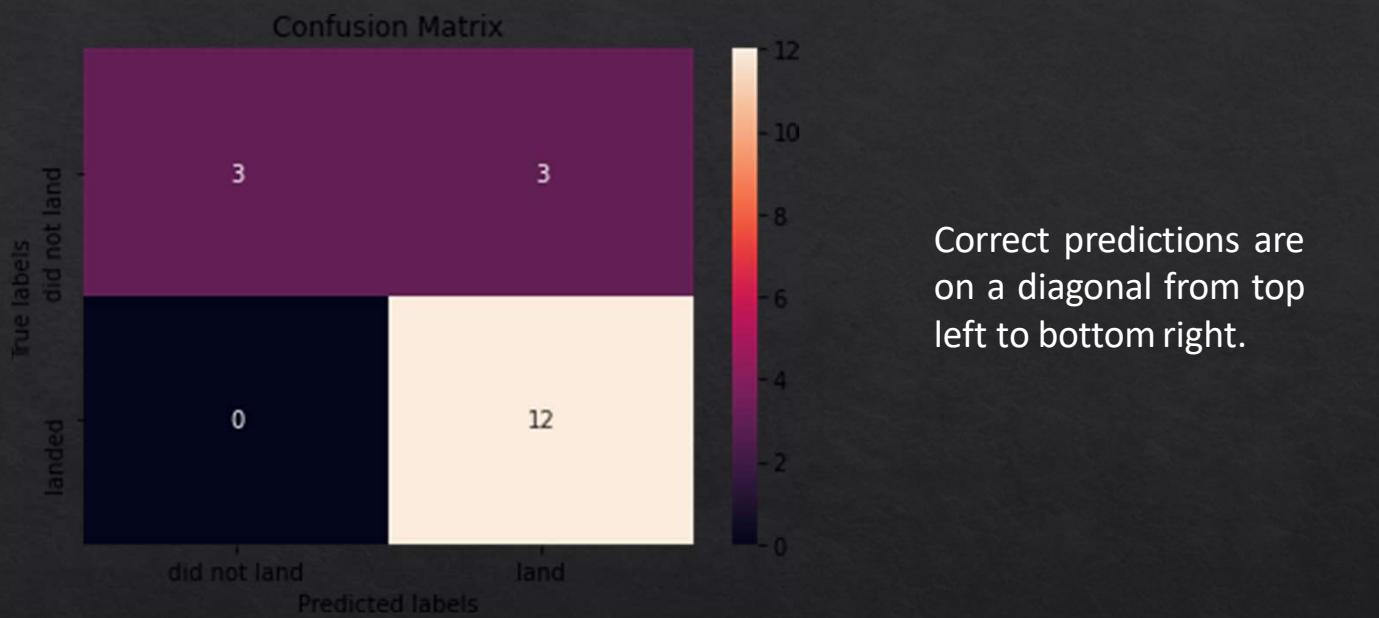


# Classification Accuracy



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.

# 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).

# CONCLUSION

---

- Our task: to develop a machine learning model for Space Y who wants to bid against SpaceX
- The goal of model is to predict when Stage 1 will successfully land to save ~\$100 million USD
- Used data from a public SpaceX API and web scraping SpaceX Wikipedia page
- Created data labels and stored data into a DB2 SQL database
- Created a dashboard for visualization
- We created a machine learning model with an accuracy of 83%
- Allon Mask of SpaceY 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 repository url:

[https://github.com/Arkadeepto2002/IBM-Data-Science-Professional-Certificate/tree/main/Applied Data Science Capstone](https://github.com/Arkadeepto2002/IBM-Data-Science-Professional-Certificate/tree/main/Applied%20Data%20Science%20Capstone)

Instructors:

**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>