

# Data Science Capstone Project

---

Zeyad Abdelmageid

<https://github.com/zeyadmageid>

13/09/2023

# Outline

---

- Executive Summary (3)
- Introduction (4)
- Methodology (6)
- Results (16)
- Conclusion (46)

# Executive Summary

---

- Gathered data from both the public SpaceX API and the SpaceX Wikipedia page. To facilitate the classification of successful landings, I established a 'class' column. I conducted data exploration through various means, such as **SQL queries, visualizations, Folium maps, and dashboards**. To prepare the data for machine learning, I selected pertinent columns to serve as features and transformed categorical variables into binary format using one-hot encoding. Following that, I standardized the data and employed **GridSearchCV** to identify the optimal parameters for our machine learning models. To assess the performance of these models, I created **visualizations** to display their accuracy scores.
- Generated four distinct machine learning models, namely ***Logistic Regression, Support Vector Machine, Decision Tree Classifier, and K Nearest Neighbors***. Each of these models exhibited comparable outcomes, achieving an accuracy rate of approximately **83.33%**. It's noteworthy that all these models tended to make an excess of predictions indicating successful landings. To enhance model precision and effectiveness, additional data is essential for improved determination and accuracy.

# Introduction

## Background:

- Space Y hired me
- Space X has best pricing (\$62 million vs. \$165 million USD)
- Largely due to the ability of recovering Stage 1 of rocket
- 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

# Methodology

---

- Data collection methodology:
  - Collected data from both SpaceX public API and SpaceX Wikipedia page
- Data wrangling:
  - Classifying true landings as successful and unsuccessful otherwise
- Perform exploratory data analysis (EDA) using visualization packages and SQL
- Perform interactive visual analytics using Folium and
- Built a web dash using plotly and dash.
- Perform predictive analysis using classification models
  - Tuned the models using GridSearchCV

# Section 1: Methodology

# Data Collection Overview

---

Data collection process involved a combination of:

- API requests from Space X public API.
- Web scraping data from Wikipedia page of Space X.

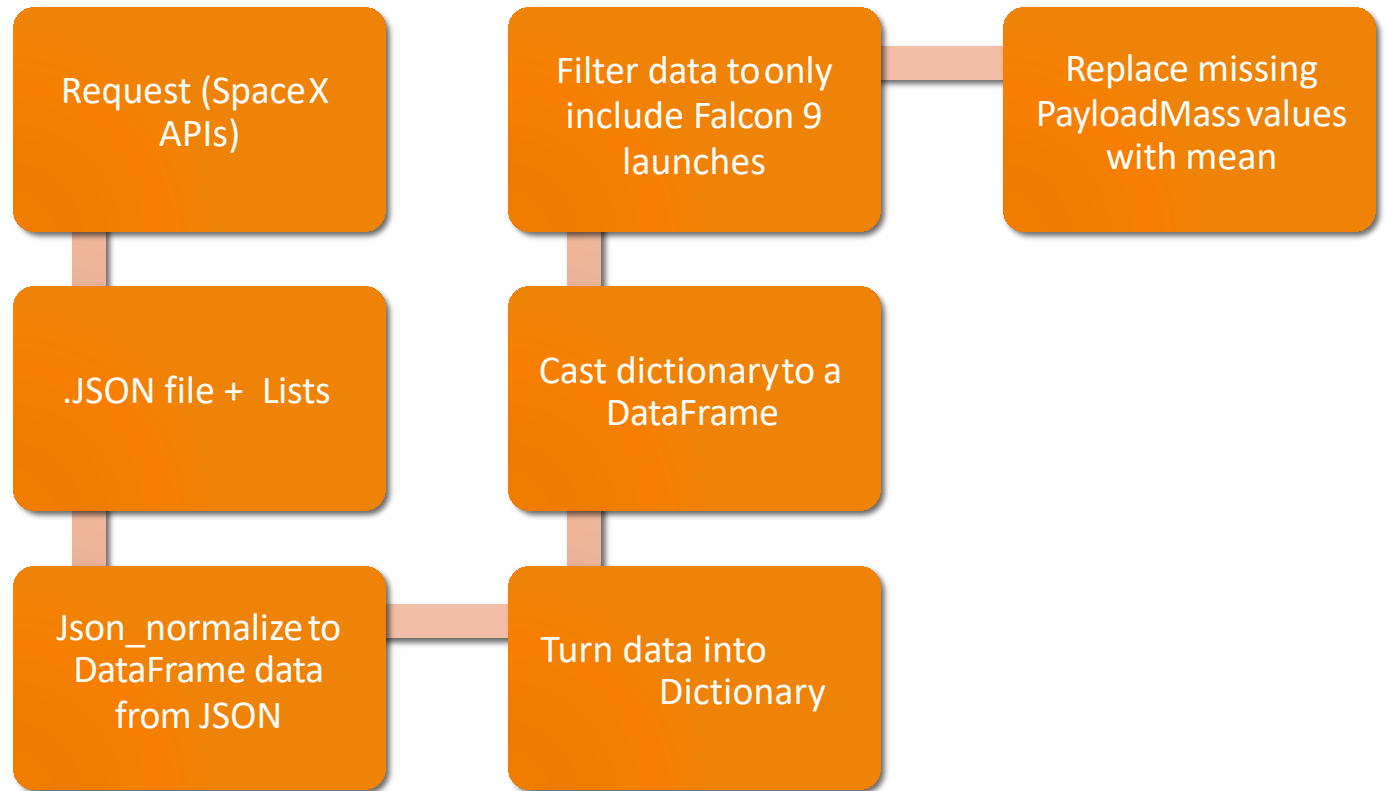
## **Space X API Data Columns:**

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

## **Wikipedia Webscrape Data Columns:**

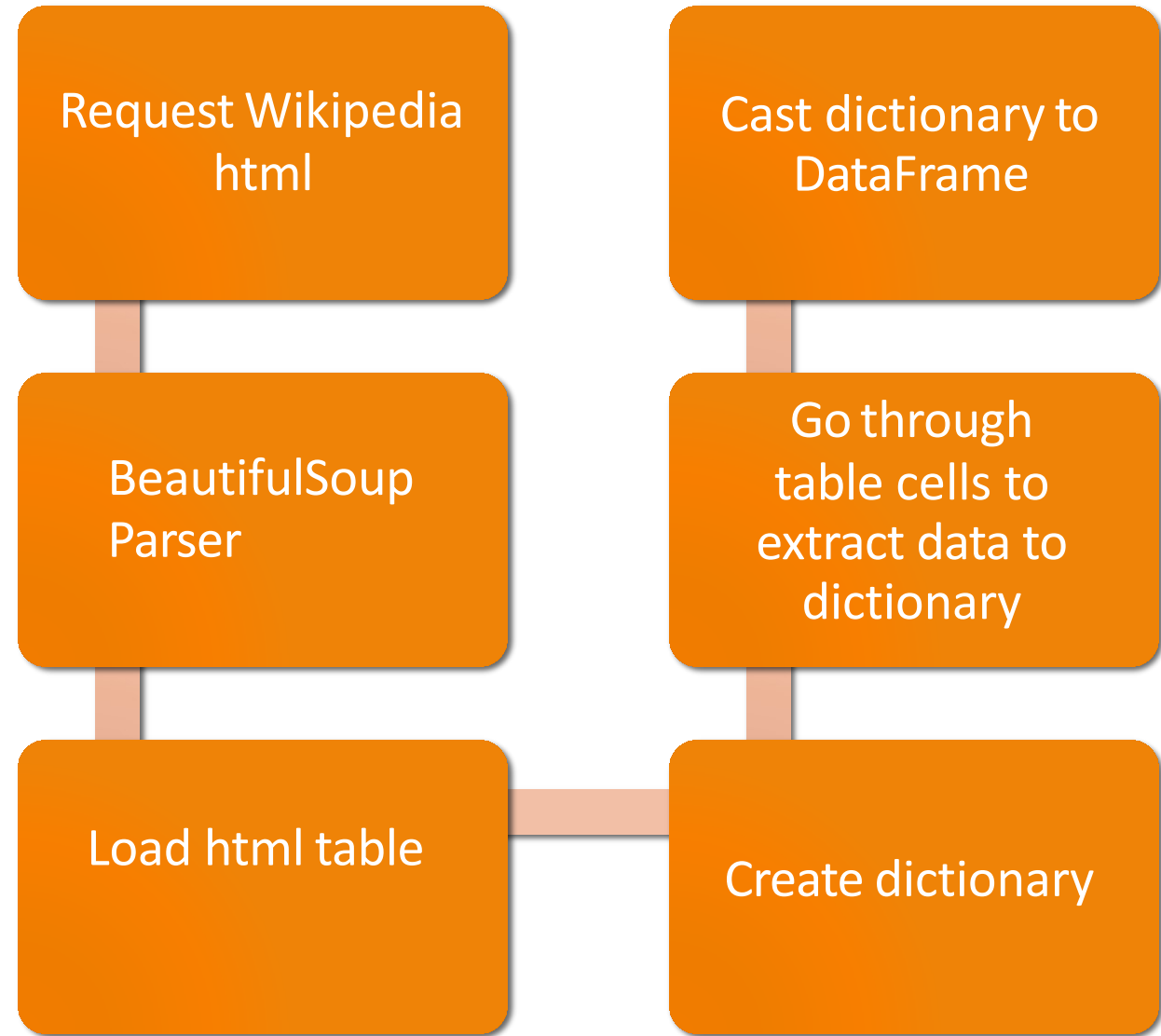
*Flight No., Launch site, Payload, PayloadMass, Orbit, Customer, Launch outcome, Version  
Booster, Booster landing, Date, Time*

# Data Collection— SpaceXAPI





# Data Collection— Web Scraping



# 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.
- 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/zeyadmageid/Capstone\\_Data\\_Science\\_Project/blob/main/3\\_Space\\_X\\_Data\\_Wrangling\\_spacex.ipynb](https://github.com/zeyadmageid/Capstone_Data_Science_Project/blob/main/3_Space_X_Data_Wrangling_spacex.ipynb)

# EDA with Data Visualization

---

EDA 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

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

## GitHub url:

[https://github.com/zeyadmageid/Capstone\\_Data\\_Science\\_Project/blob/main/5\\_Space\\_X\\_EDA\\_DataViz\\_Using\\_Pandas\\_and\\_Matplotlib\\_SpaceX.ipynb](https://github.com/zeyadmageid/Capstone_Data_Science_Project/blob/main/5_Space_X_EDA_DataViz_Using_Pandas_and_Matplotlib_SpaceX.ipynb)

# EDA with SQL

---

Loaded data set into IBM DB2 Database.

Queried using SQL magic function.

Queried information about launch site names, mission outcomes, various payload sizes of customers and booster versions, and landing outcomes

GitHub url:

[https://github.com/zeyadmageid/Capstone\\_Data\\_Science\\_Project/blob/main/4\\_Space\\_X\\_EDA\\_Using\\_SQL.ipynb](https://github.com/zeyadmageid/Capstone_Data_Science_Project/blob/main/4_Space_X_EDA_Using_SQL.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/zeyadmageid/Capstone\\_Data\\_Science\\_Project/blob/main/6\\_Space\\_X\\_Launch\\_Sites\\_Locations\\_Analysis\\_with\\_Folium\\_Interactive\\_Visual\\_Analytics.ipynb](https://github.com/zeyadmageid/Capstone_Data_Science_Project/blob/main/6_Space_X_Launch_Sites_Locations_Analysis_with_Folium_Interactive_Visual_Analytics.ipynb)

# 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:

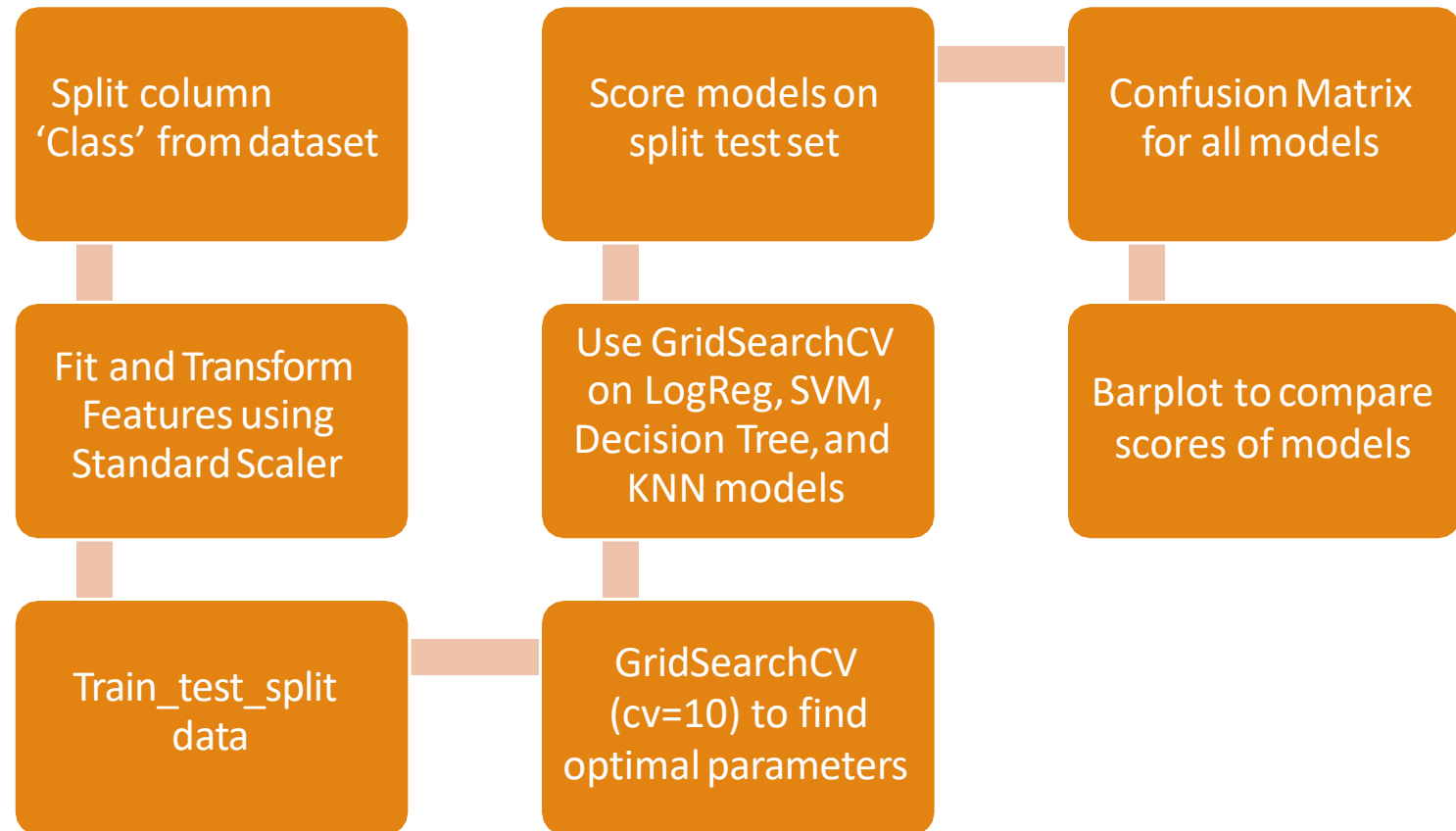
[https://github.com/zeyadmageid/Capstone\\_Data\\_Science\\_Project/blob/main/7.%20Build%20an%20Interactive%20Dashboard%20with%20Plotly%20Dash%20-%20spacex\\_dash\\_app.py](https://github.com/zeyadmageid/Capstone_Data_Science_Project/blob/main/7.%20Build%20an%20Interactive%20Dashboard%20with%20Plotly%20Dash%20-%20spacex_dash_app.py)

# Predictive analysis (Classification)

---

GitHub url:

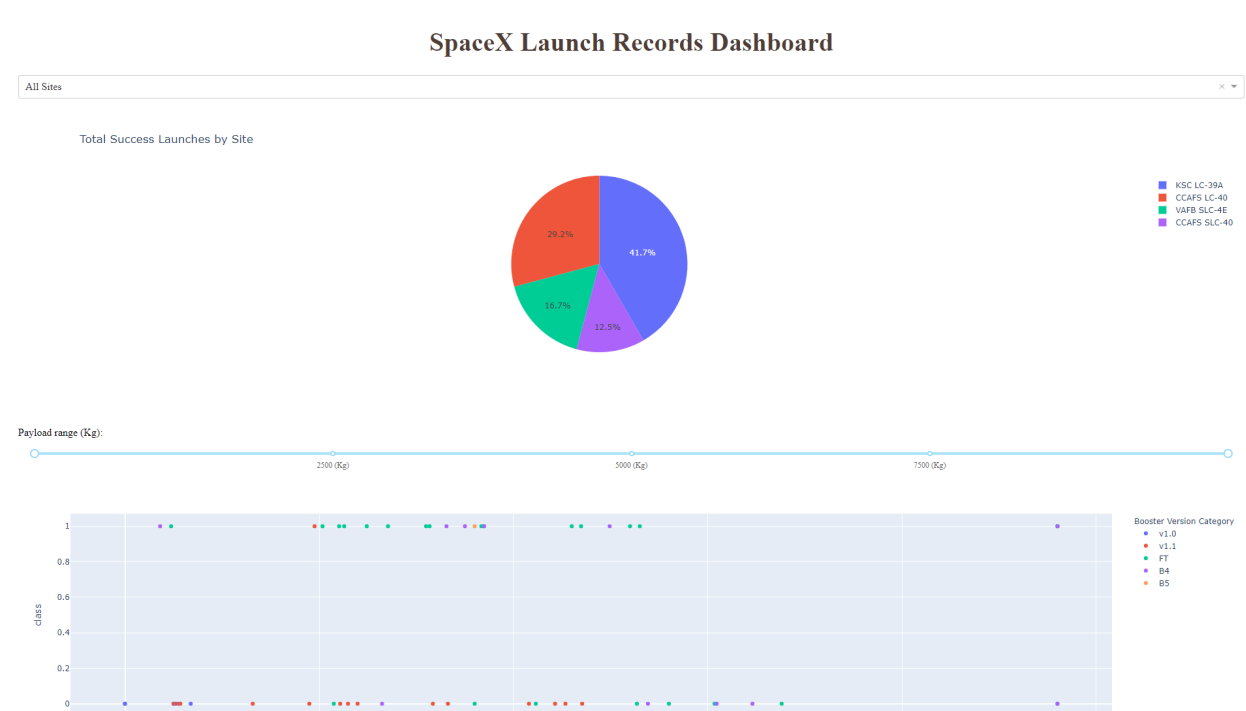
[https://github.com/zeyadmaheid/Capstone\\_Data\\_Science\\_Project/blob/main/8\\_Space\\_X\\_Machine\\_Learning\\_Prediction.ipynb](https://github.com/zeyadmaheid/Capstone_Data_Science_Project/blob/main/8_Space_X_Machine_Learning_Prediction.ipynb)



# Section 2: Results



# Results

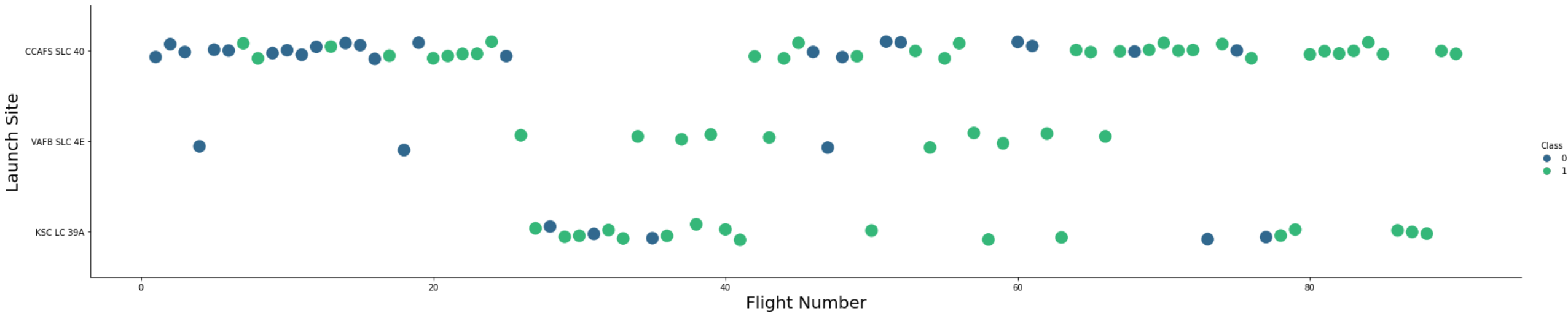


Plotly dashboard

# EDA with Visualization

EXPLORATORY DATA ANALYSIS

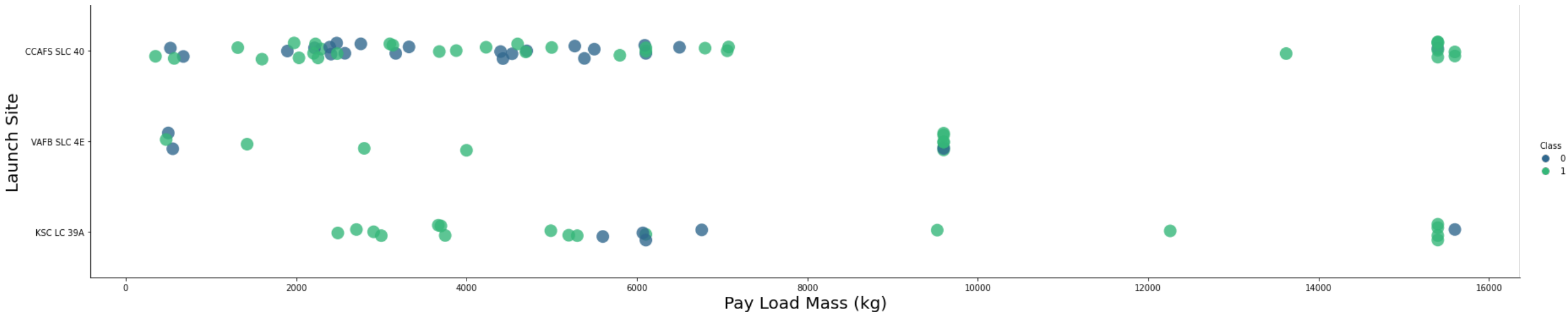
# Flight Number vs. LaunchSite



Green indicates successful launch; Purple indicates unsuccessful launch.

Graphic suggests an increase in success rate as more flights are deployed (Flight Number gets higher). CCAFS appears to be the main launch site as it has the most volume.

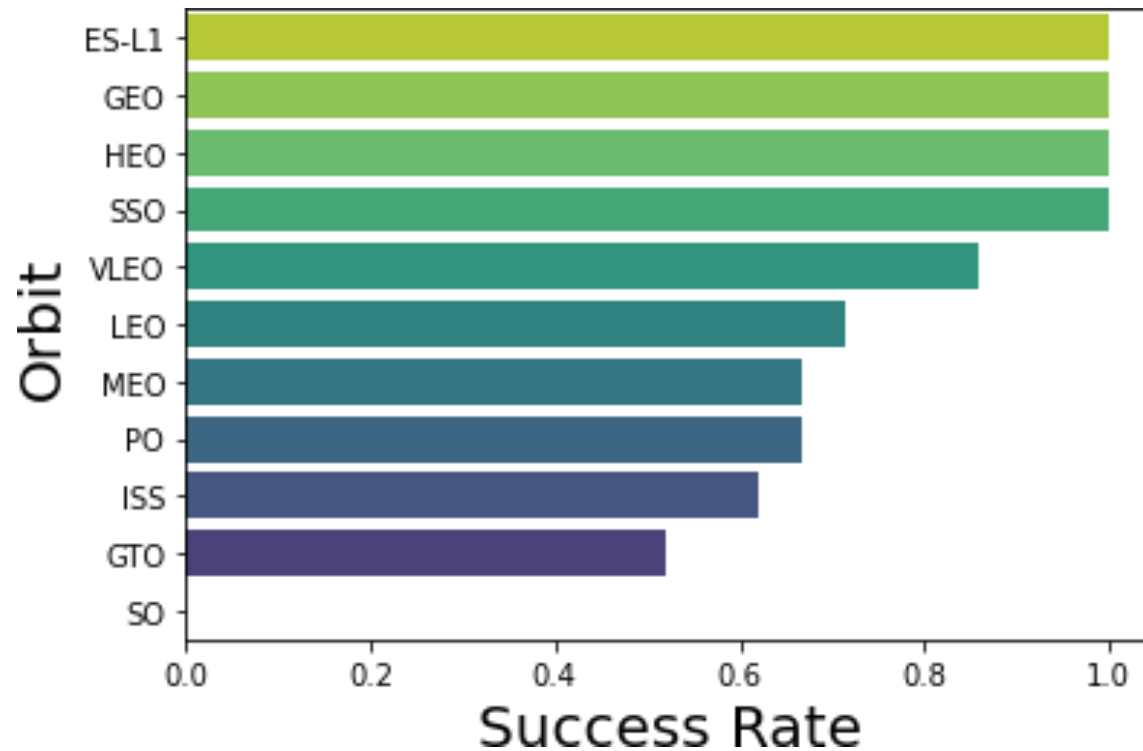
# Payload vs. Launch Site



Green indicates successful launch; Purple indicates unsuccessful launch.

Payload mass appears to fall mostly between 0-7000 kg.

# Successrate vs. Orbittype



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

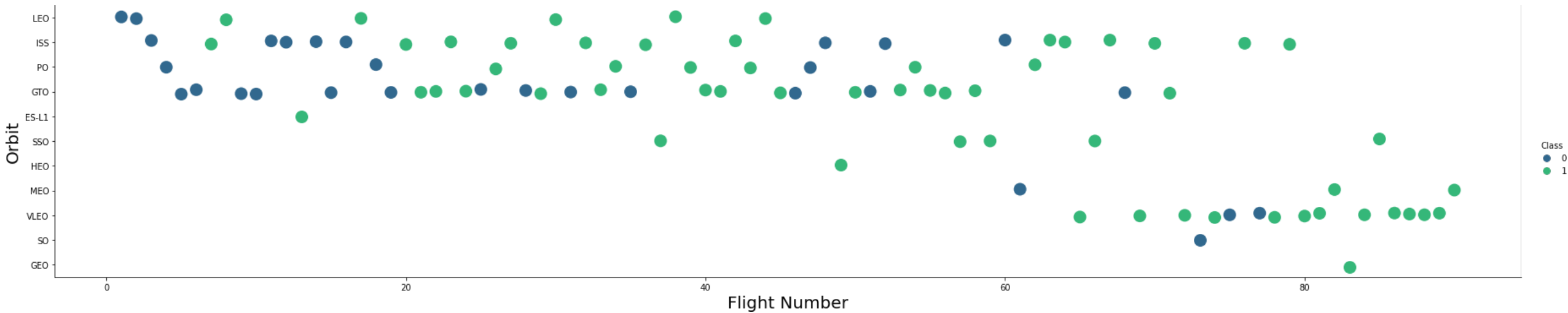
SSO (5) has 100% success rate

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



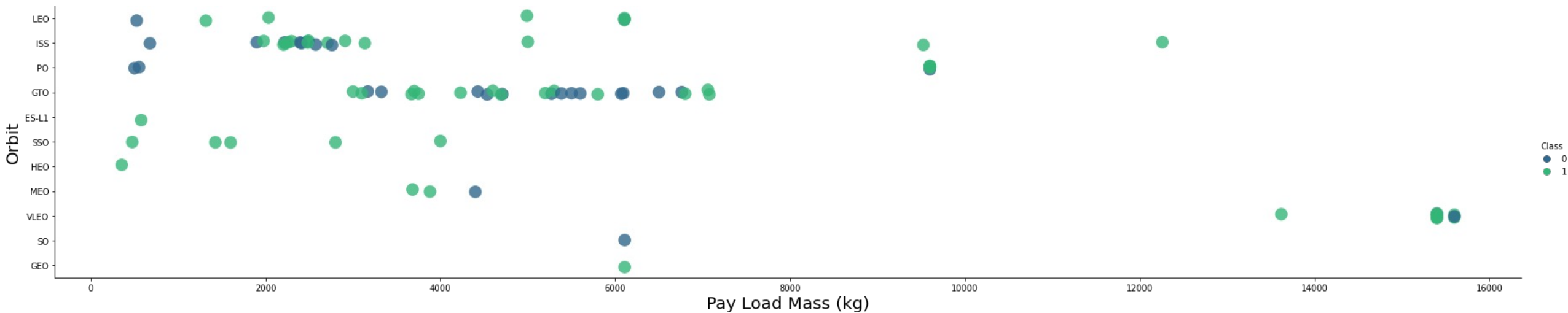
Green indicates successful launch; Purple indicates unsuccessful launch.

Launch Orbit changed over Flight Number.

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

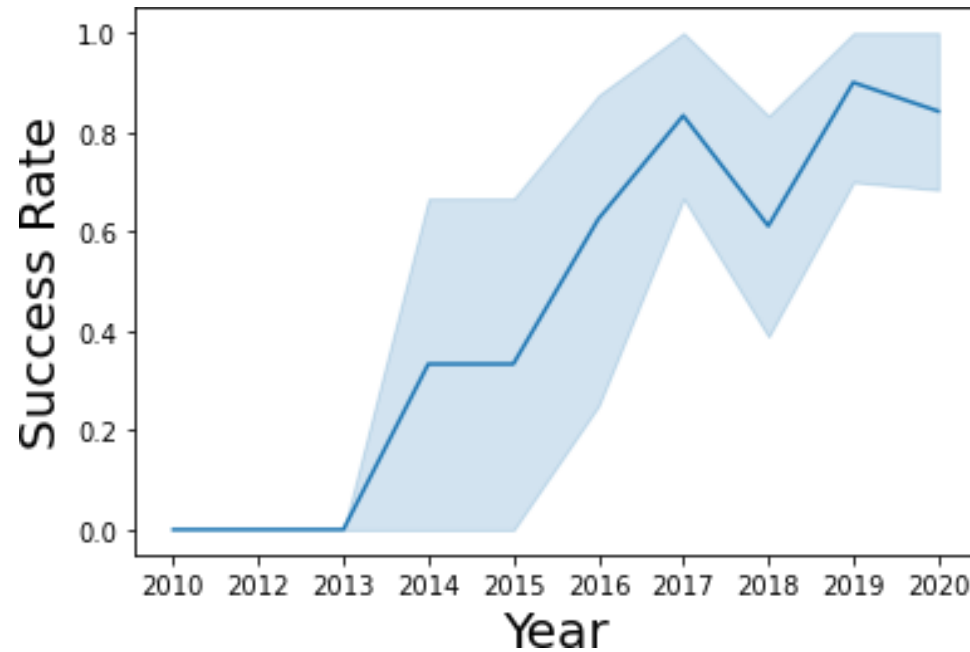


Green indicates successful launch; Purple indicates unsuccessful launch.

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



95% confidence interval  
(light blue shading)

Success generally increases over time since 2013 with a slight dip starting in 2017 and back up again in 2019

Success in recent years at around 80-90%



# EDA with SQL

# All Launch Site Names

---

```
In [4]: %%sql
        SELECT UNIQUE LAUNCH_SITE
        FROM SPACEXDATASET;

* ibm_db_sa://ftb12020:***@0c77d6f:
Done.
```

Query unique launch site names from database.

```
Out[4]:
```

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

# Launch Site Names 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

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 when NASA was the customer.

# 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:***@0c77d6f2-5da9-48a9-81f8-86
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 v1.1 is 2928 kgs

# 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 was on the end of 2015.

# 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

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

# 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

```
%%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 booster version.

# 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.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.

# 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.bs2io90l08kqb1od8lce
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

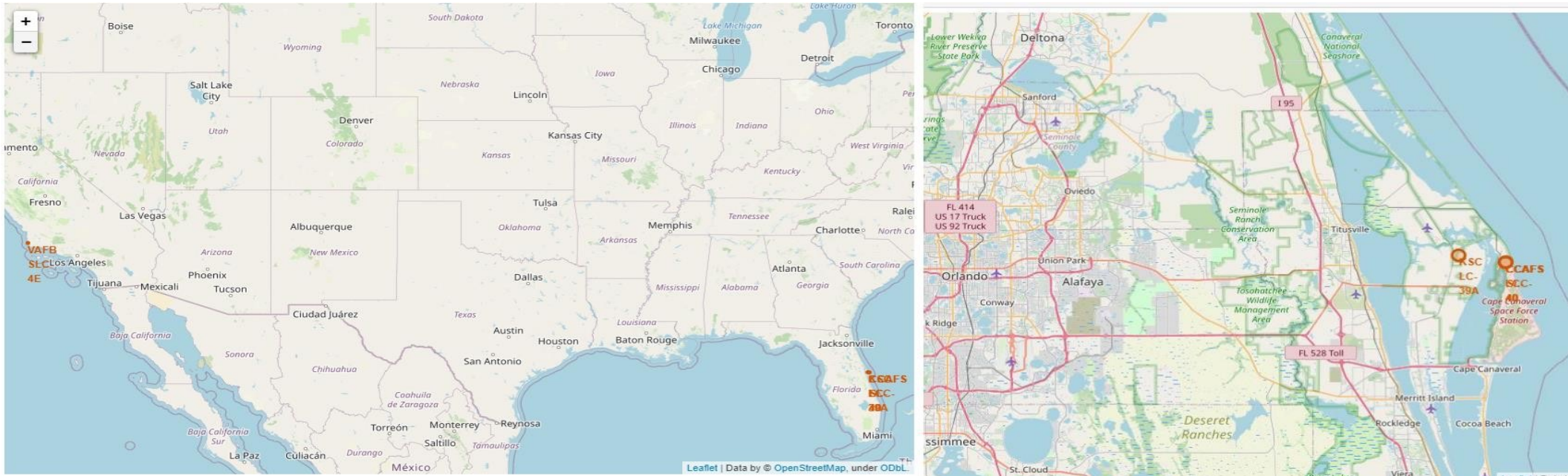
There are two types of successful landing outcomes: drone ship and ground pad landings.

There were 5 successful drone landings and 3 on ground pad in total during this time period

# Interactive Map with Folium

# Launch Site

## Locations

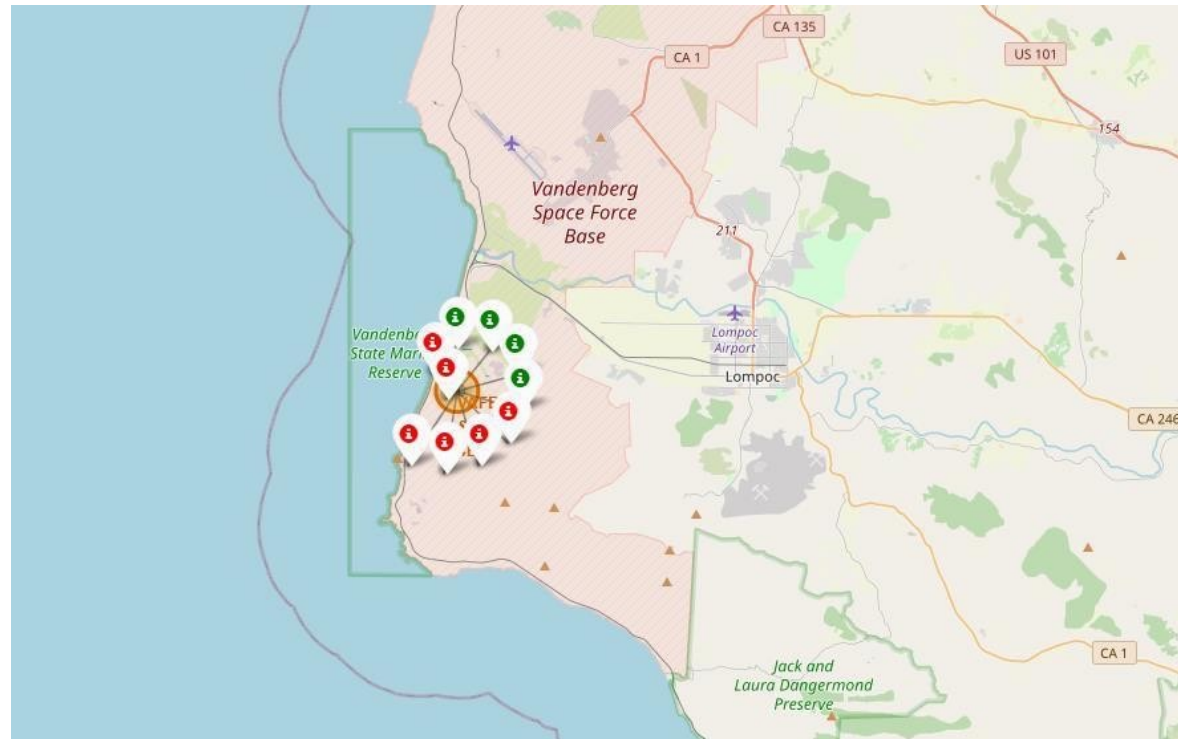


The left map shows all launch sites relative US map.

The right map shows the two Florida launch sites since they are very close to each other.

# Color-Coded Launch Markers

---

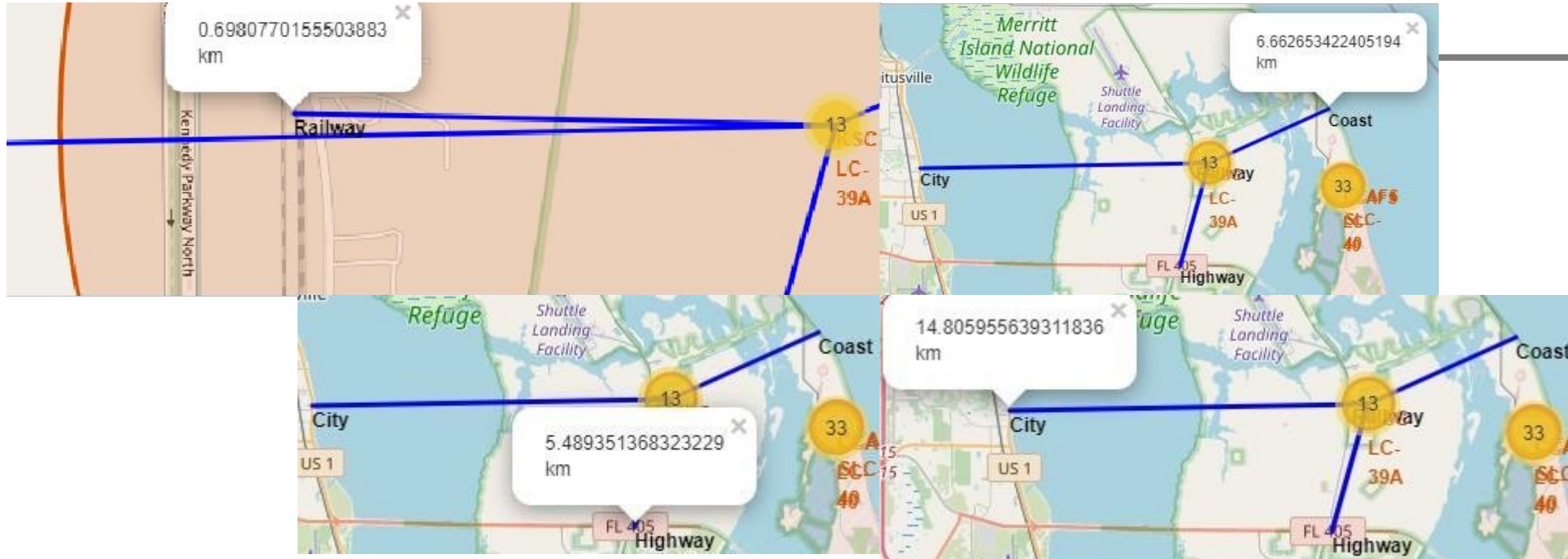


Clusters on Folium map can be clicked on to display each successful landing (green icon) and failed landing (red icon).



# Key Location

## Proximities



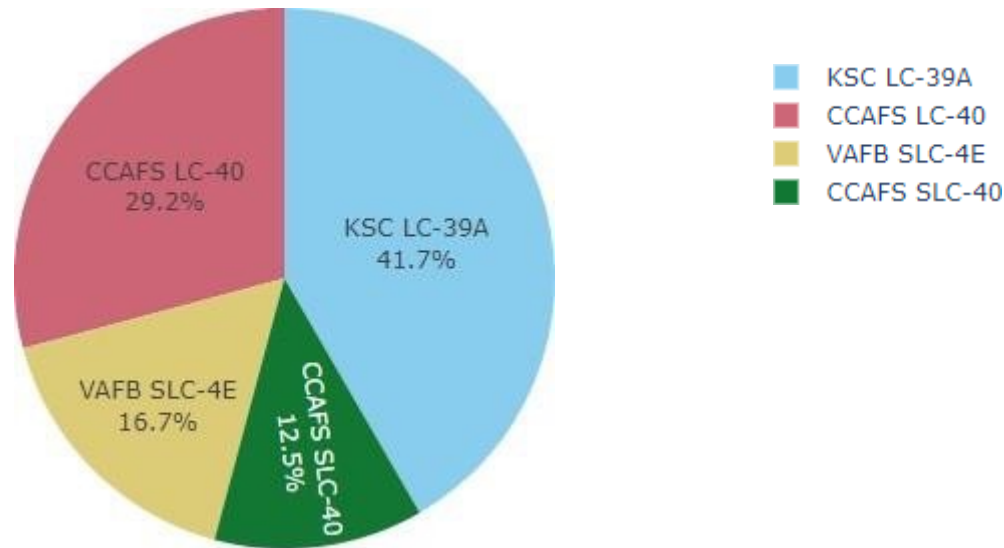
Launch sites are close to highways for human and supply transport. Launch sites are also close to coasts and relatively far from cities so that launch failures can land in the sea to avoid rockets falling on populated areas.

# Build a Dashboard with Plotly Dash



# Successful Launches Across Launch Sites

---

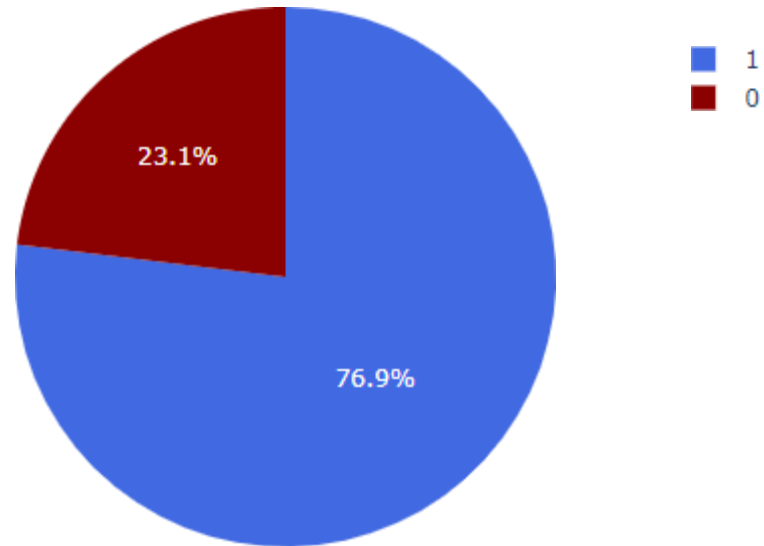


This is the distribution of successful landings across all launch sites. CCAFS LC-40 is the old name of CCAFS SLC-40 so CCAFS and KSC have the same amount of successful landings, but a majority of the successful landings were performed before the name change. VAFB has the smallest share of successful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

Payload range (Kg):



Payload Mass vs. Success vs. Booster Version Category



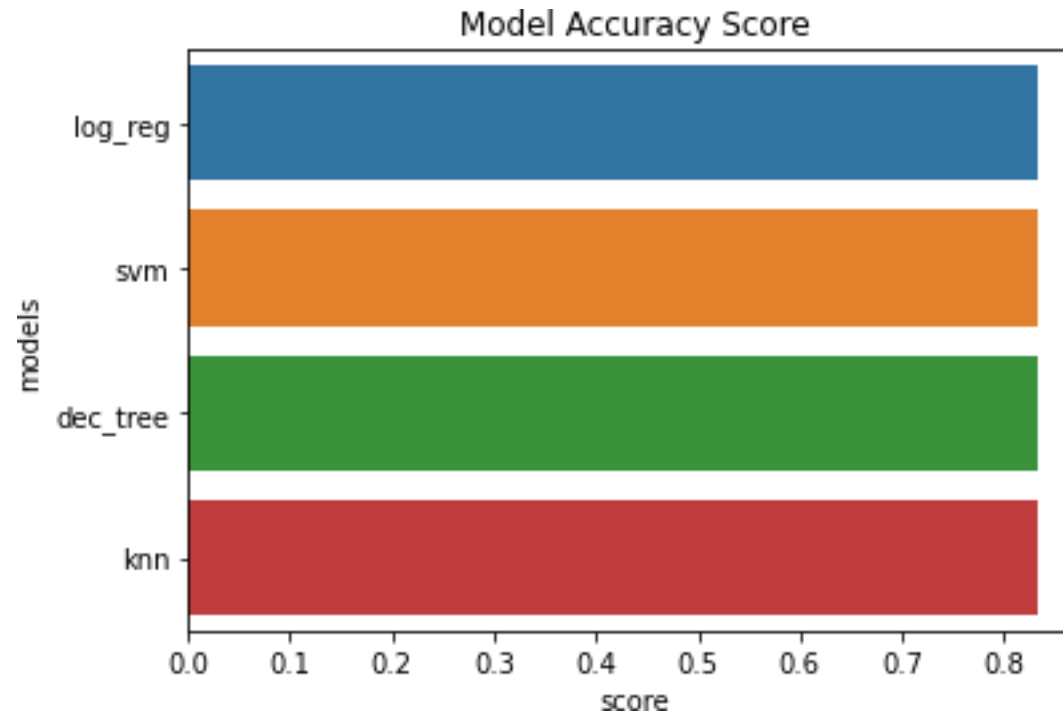
Plotly dashboard has a Payload range selector. However, this is set from 0-10000 instead of the max Payload of 15600. Class indicates 1 for successful landing and 0 for failure. Scatter plot also accounts for booster version category in color and number of launches in point size. In this particular range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

- PredictiveAnalysis(Classification)

---

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

# Classification Accuracy



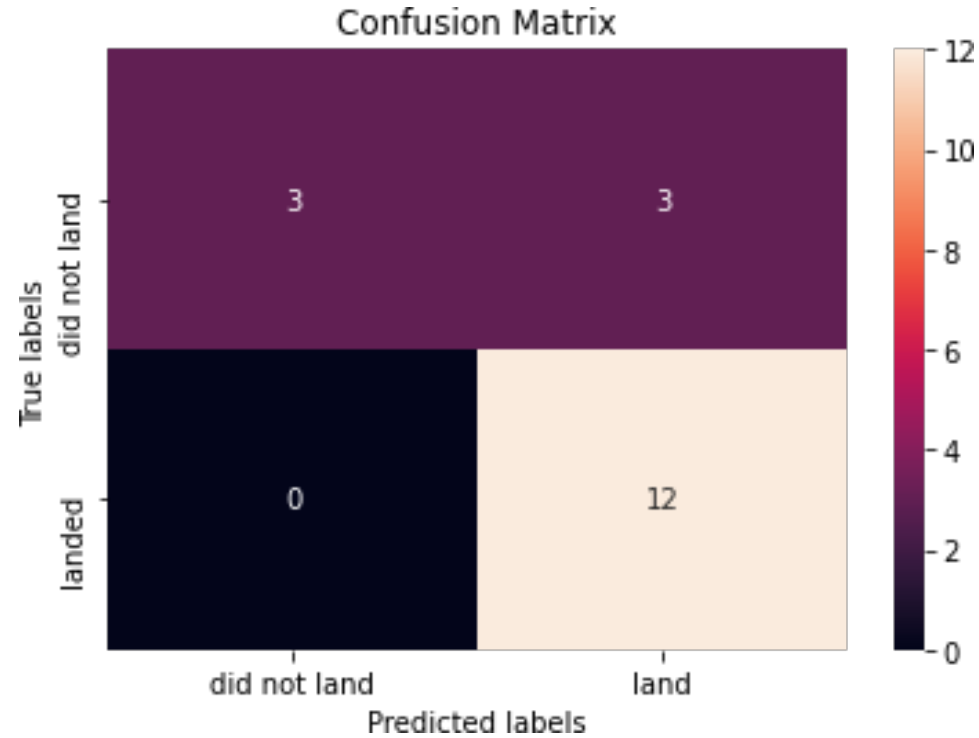
All models had virtually the same accuracy on the test set at 83.33% accuracy.

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.

# Section 3: Conclusion

# CONCLUSION

---

- Our task: to develop a machine learning model for SpaceY who wants to bid against SpaceX
- The goal of model is to predict when Stage 1 will successfully land
- Used data from a public SpaceX API and web scraping SpaceX Wikipedia page
- Created data labels and stored data into a SQL database
- Created a dashboard for visualization
- Created a machine learning model with an accuracy of 83% after tuning