# Data Science Capstone Project

Aly Saleh

https://github.com/AlyHSaleh/IBMCapstoneProject

22/03/2024

# Outline

- Executive Summary
- Introduction
- Methodology
- Results
- Conclusion

# 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 best parameters for machine learning models. Visualize accuracy score of all models.

- Four machine learning models were produced: Logistic Regression, Support Vector Machine, Decision Tree Classifier, and K Nearest Neighbors. All produced similar results with accuracy rate of about 83.33%. All models over predicted successful landings. 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 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

# 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

# Methodology

OVERVIEW OF DATA COLLECTION, WRANGLING, VISUALIZATION, DASHBOARD, AND MODEL METHODS

# 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,
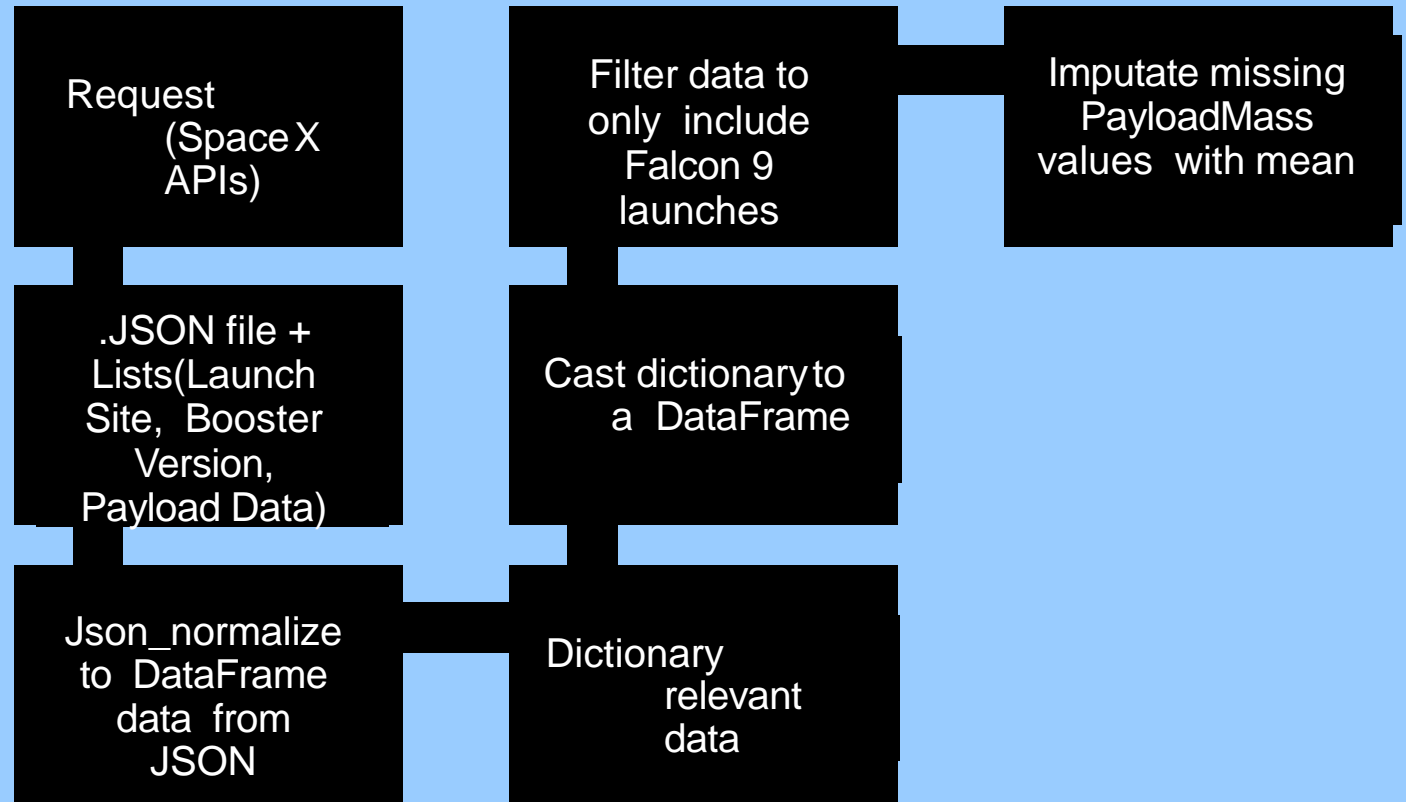Reused, Legs, LandingPad, Block, ReusedCount, Serial, Longitude, Latitude

Wikipedia Webscrape Data Columns:

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

# Data Collection – SpaceX API

Request (Space X APIs)

Filter data to only include Falcon 9 launches

Imputate missing PayloadMass values with mean

.JSON file + Lists(Launch Site, Booster Version, Payload Data)

Cast dictionary to a DataFrame

Json_normalize to DataFrame data from JSON

Dictionary relevant data

# Data Collection – Web Scraping

Request Wikipedia html

BeautifulSoup html5lib Parser

Find launch info html table

Create dictionary

Iterate through table cells to extract data to dictionary

Cast dictionary to DataFrame

# 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/AlyHSaleh/IBMCapstoneProject/blob/main/labs-jupyter-spacex-Data%20wrangling.ipynb

# EDA with SQL

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/AlyHSaleh/IBMCapstoneProject/blob/main/jupyter-labs-eda-sql-coursera_sqllite.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 to
decide if a relationship exists so that they could be used in training the machine learning model

GitHub url:

https://github.com/AlyHSaleh/IBMCapstoneProject/blob/main/jupyter-labs-eda-dataviz.ipynb.jupyterlite.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/AlyHSaleh/IBMCapstoneProject/blob/main/lab_jupyter_launch_site_location.jupyterlite.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 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.
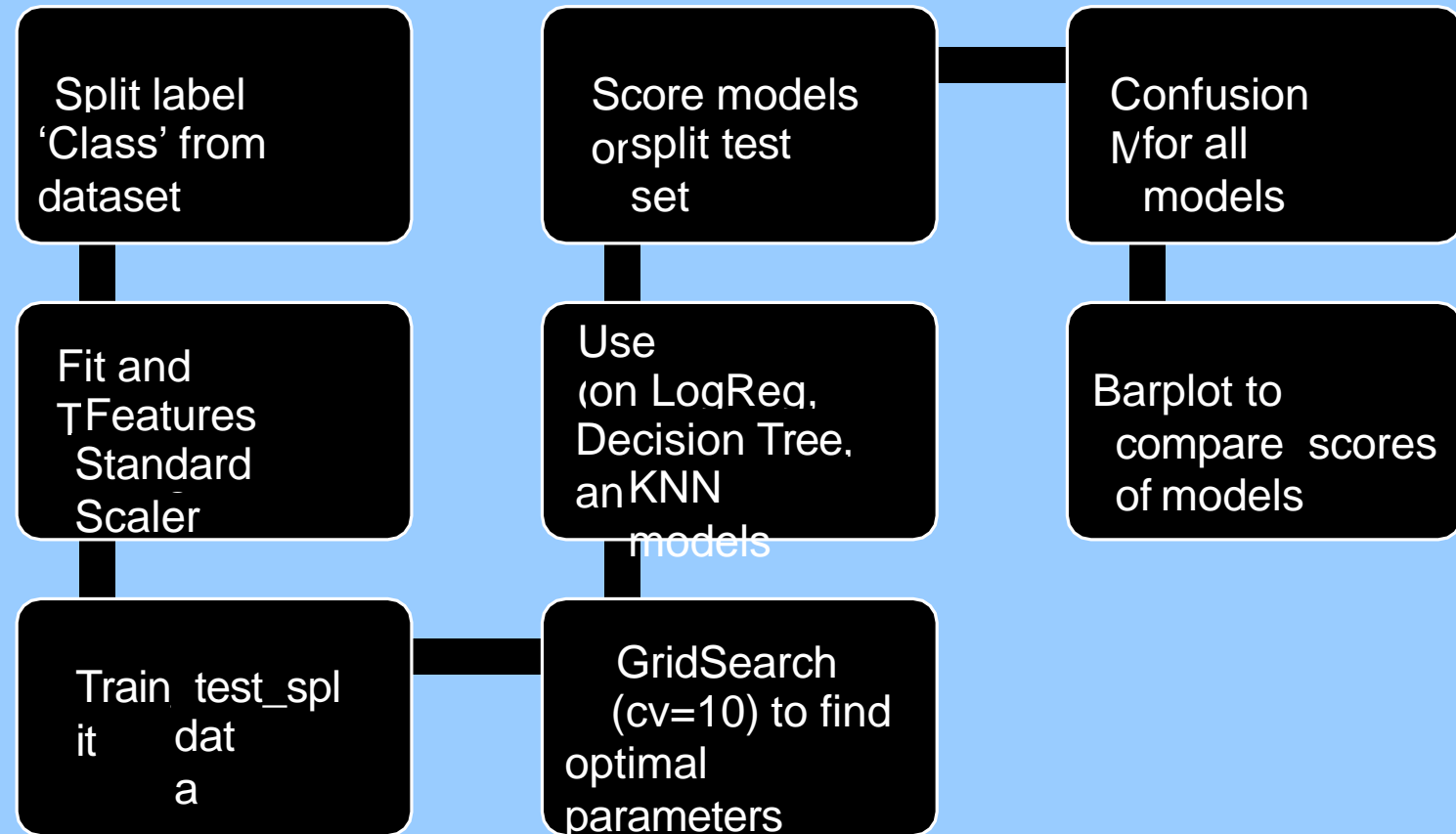
GitHub url:

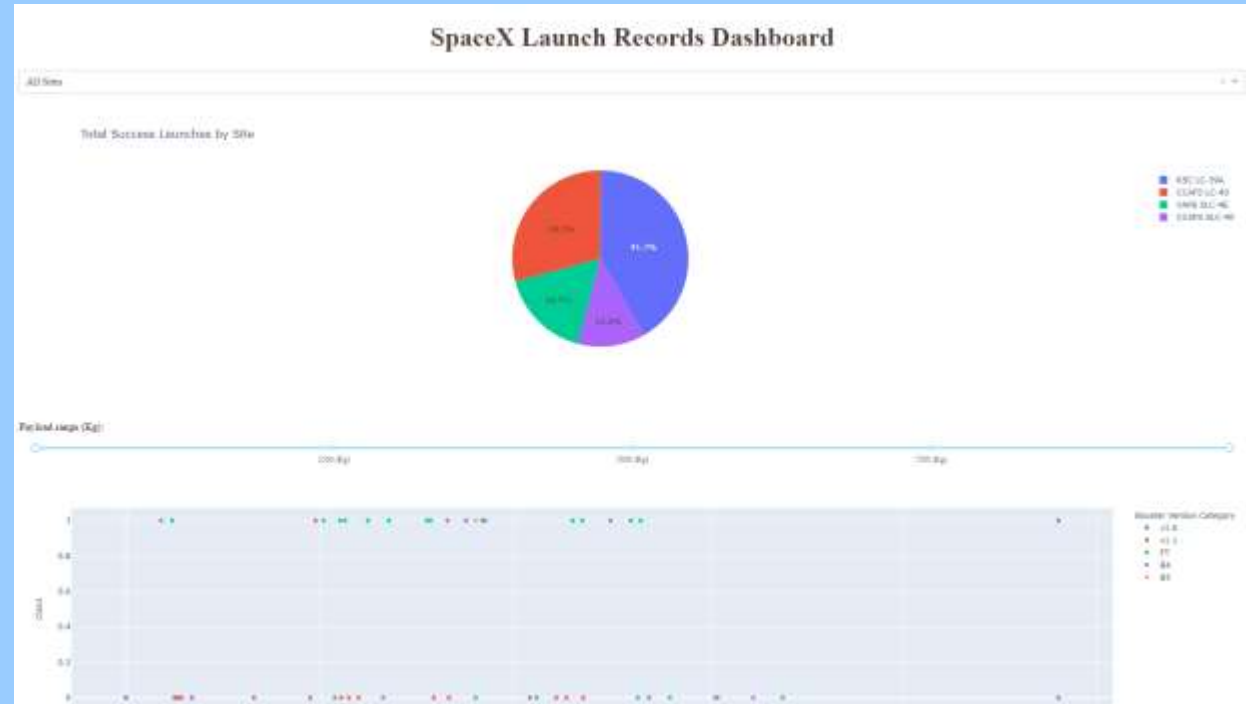https://github.com/AlyHSaleh/IBMCapstoneProject/blob/main/spacex_dash_app.py

# Predictive analysis (Classification)

GitHub url:
https://github.com/AlyHSal
eh/IBMCapstoneProject/bl
ob/main/SpaceX_Machine
_Learning_Prediction_Part
_5.jupyterlite.ipynb

Split label 'Class' from dataset

Fit and Features Standard Scaler

Train test_split data

Score models or split test set

Use (on LogReg, Decision Tree, an KNN models

GridSearch (cv=10) to find optimal parameters

Confusion M for all models
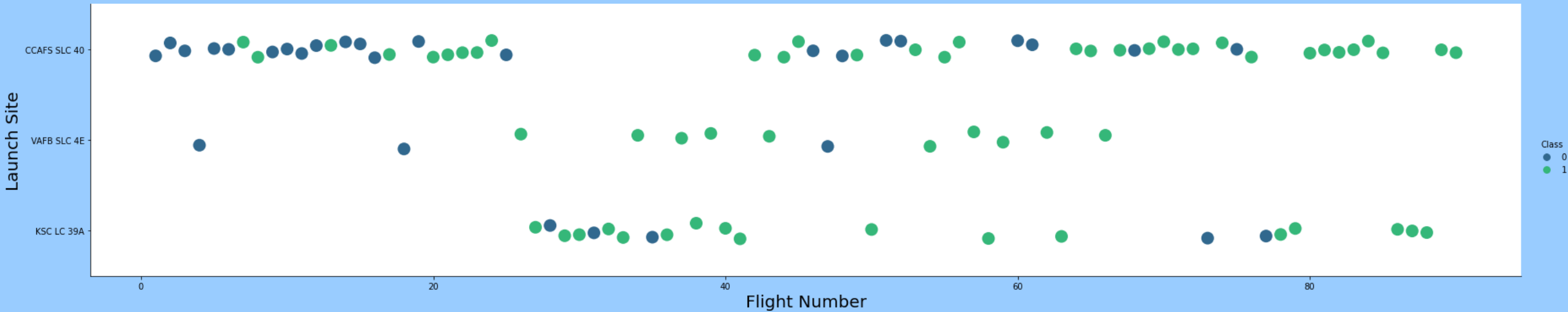
Barplot to compare scores of models

# Results



This is a preview of the Plotly dashboard. The following sides 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.

# E D A with Visualization
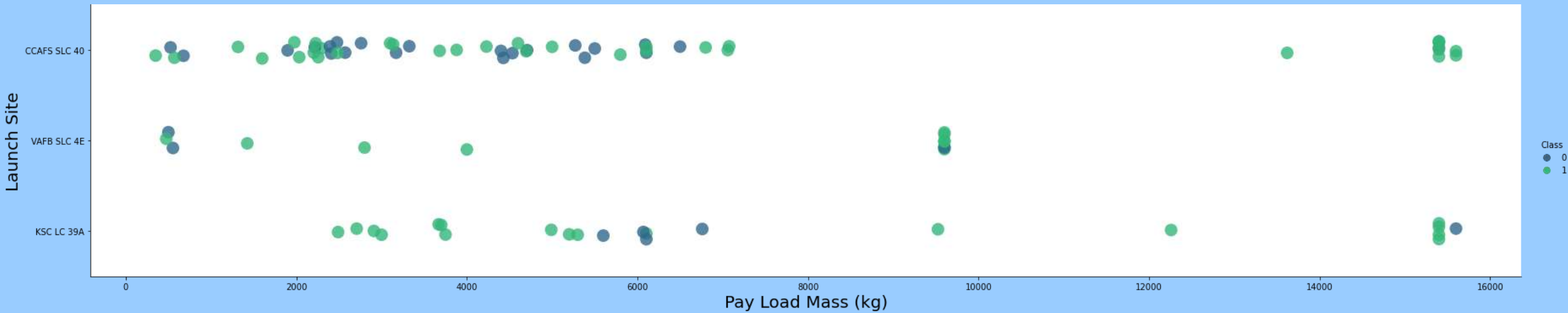
EXPLORATORY  DATA ANALYSIS   WITH  SEABORN  PLOTS

# 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.
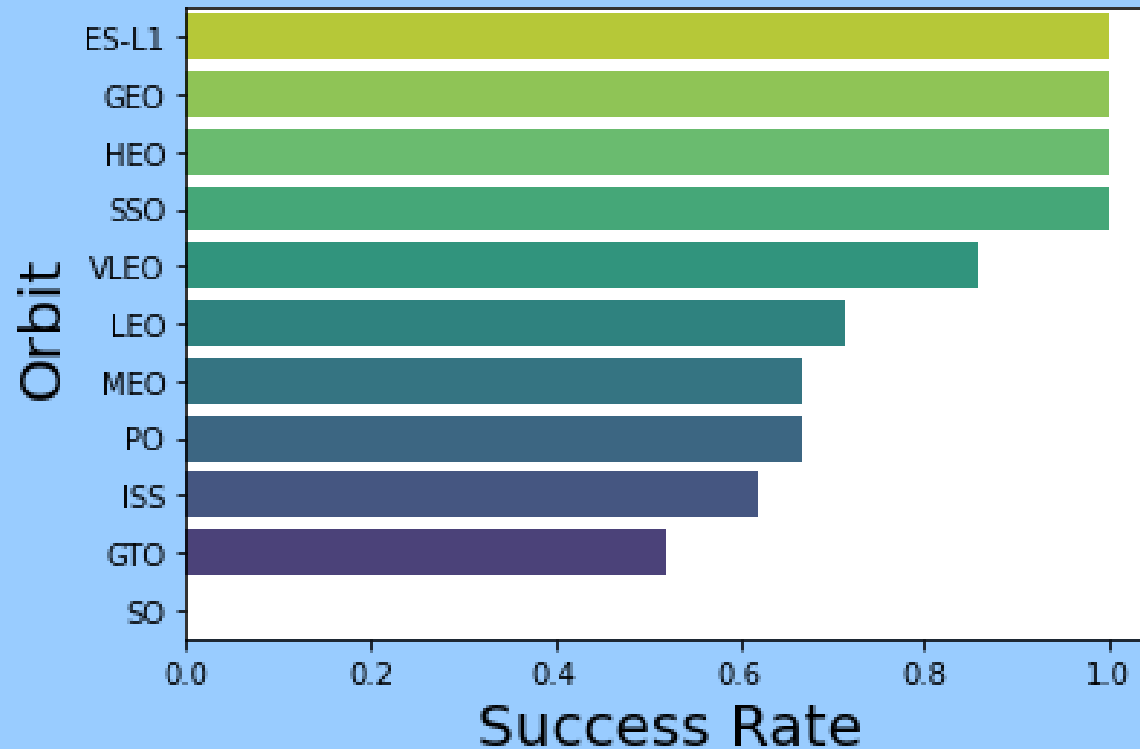
# Payload vs. Launch Site



Green indicates successful launch; Purple indicates unsuccessful launch.

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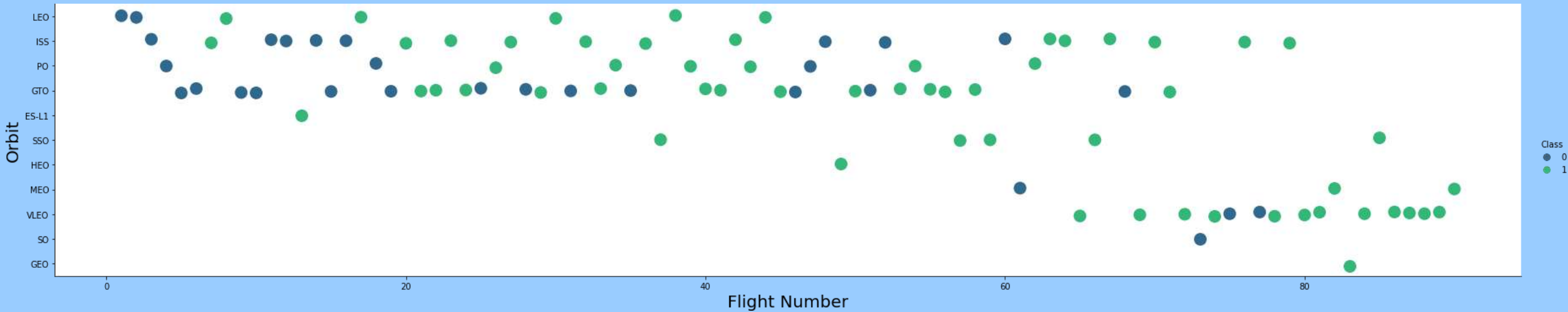Different launch sites also seem to use different payload mass.

# Success rate vs. Orbit type



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

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
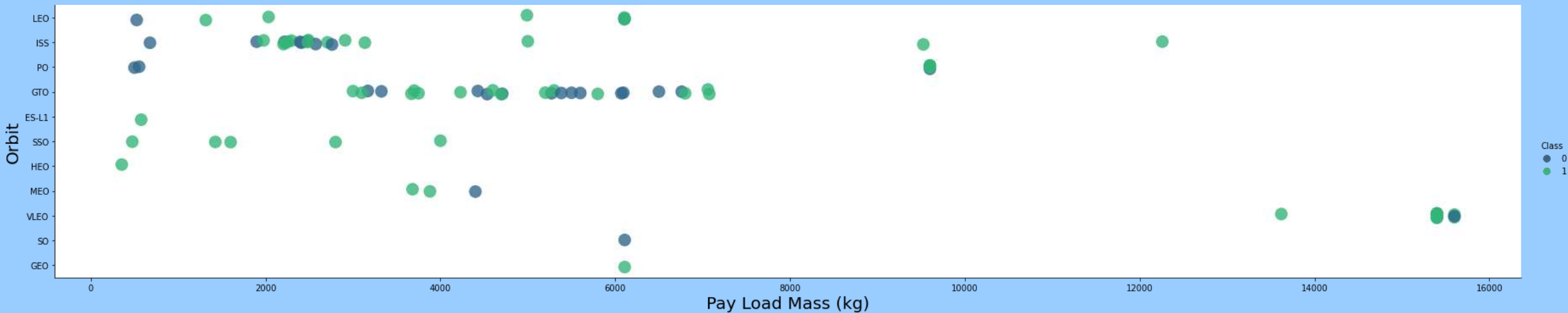
# Flight Number vs. Orbittype



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
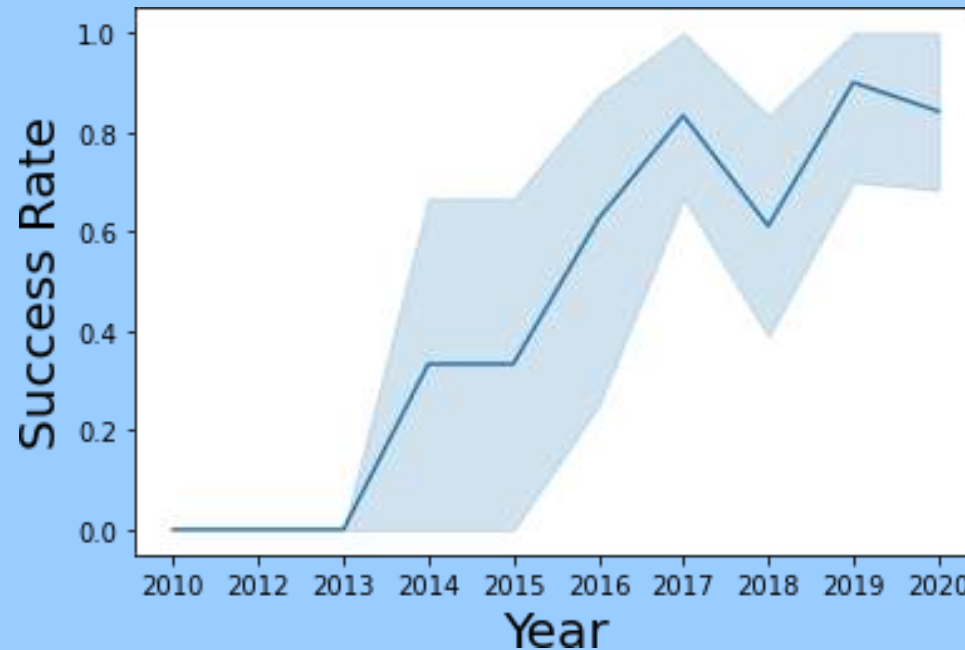
# Payload vs. Orbit type



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



95% confidence interval  (light blue shading)

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

Success in recent years at around 80%

# EDAwith SQL

EXPLORATORY  DATA ANALYSIS   WITH  SQL  DB2

INTEGRATED  IN  PYTHON  WITH  SQLALCHEMY

# All Launch Site Names

```
%%sql
select DISTINCT LAUNCH_SITE from SPACEXTBL
```

* sqlite:///my_data1.db
Done.
Out[10]:

**Launch_Site**

CCAFS LC-40

VAFB SLC-4E

KSC LC-39A

CCAFS SLC-40

Query unique launch site names from database.

Results is 4 unique launch_site

values:

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

# Launch Site Names Beginning with `CCA`

```
In [11]:    %%sql
            select *
            from SPACEXTBL
            where launch_site like 'CCA%' limit 5
```

```
 * sqlite:///my_data1.db
Done.
```

Out[11]:

| 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 | 7:44:00 | F9 v1.0 B0005 | CCAFS LC-40 | Dragon demo flight C2 | 525 | LEO (ISS) | NASA (COTS) | Success | No attempt |
| 2012-10-08 | 0:35:00 | 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

```
In [12]:   %%sql
           select sum(payload_mass__kg_) as sum
           from SPACEXTBL
           where customer like 'NASA (CRS)'

         * sqlite:///my_data1.db
         Done.

Out[12]:    sum
            _____

            45596
```

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

# Average Payload Mass by F9 v1.1



```
In [13]:   %%sql
           select avg(payload_mass__kg_) as Average
           from SPACEXTBL
           where booster_version like 'F9 v1.1%'

            * sqlite:///my_data1.db
           Done.

Out[13]:
                    Average

           2534.6666666666665
```

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

# First Successful Ground Pad Landing Date

```
In [14]:  %%sql
          select min(date) as Date
          from SPACEXTBL
          where mission_outcome like 'Success'
```

```
 * sqlite:///my_data1.db
Done.
```

Out[14]:

| Date |
| --- |
| 2010-06-04 |

This query returns the first successful ground pad landing  date.

# Successful Drone Ship Landing with Payload Between 4000 and 6000

```
In [16]:    %%sql
            select booster_version
            from SPACEXTBL
            where (mission_outcome like 'Success')
            AND (payload_mass__kg_ BETWEEN 4000 AND 6000)
            AND (landing_outcome like 'Success (drone ship)')

             * sqlite:///my_data1.db
            Done.
Out[16]:    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.

# 2015 Failed Drone Ship Landing Records



This query returns the Date, 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.
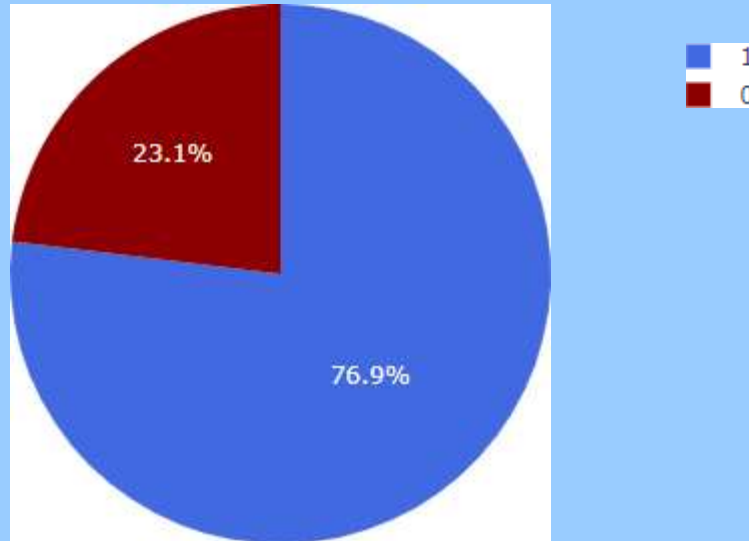
# Build a Dashboard with Plotly Dash

# 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 where 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.

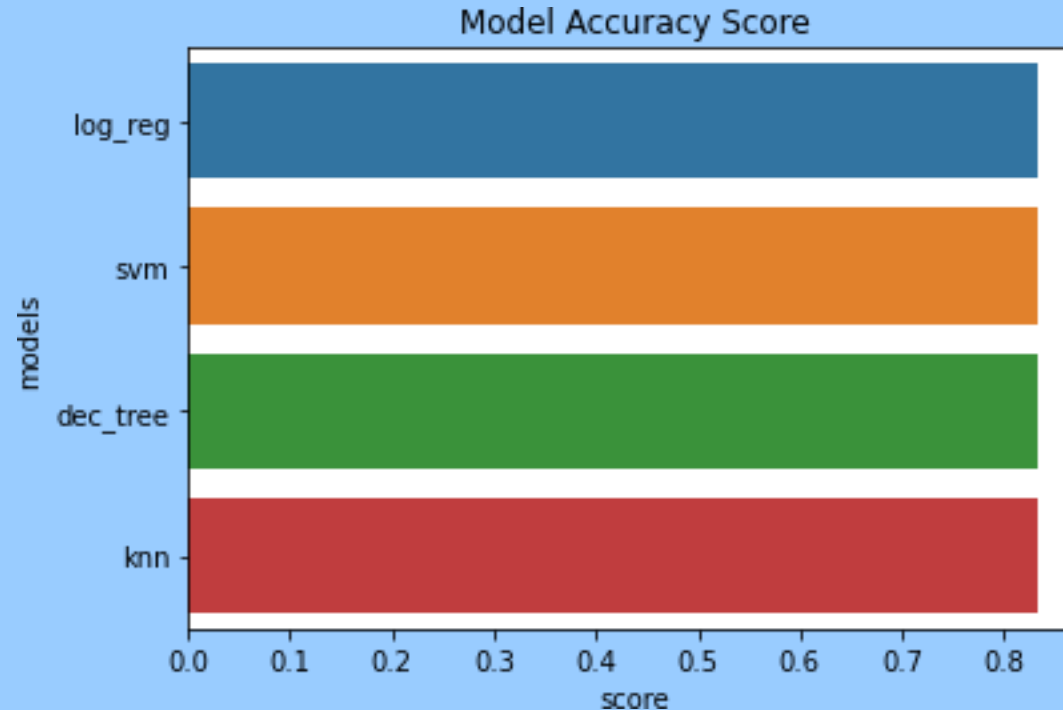# Payload Mass vs. Success vs. Booster Version Category



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.

# Predictive Analysis (Classification)

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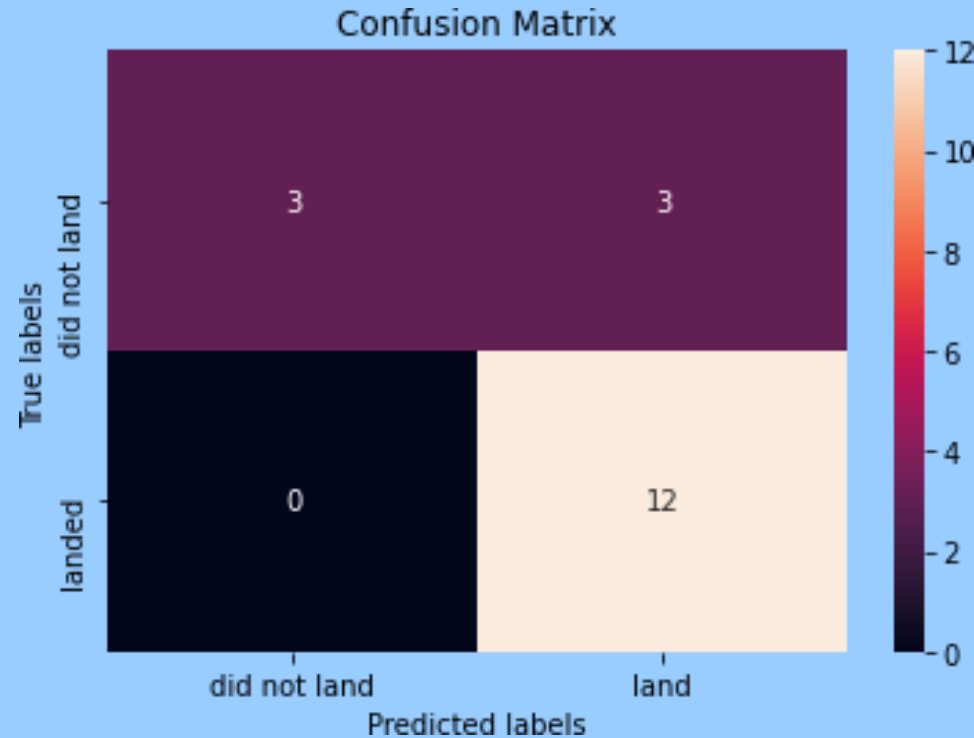
# Classification Accuracy



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



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

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.

# 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
- Used data from a public SpaceX API and web scraping SpaceX Wikipedia page
- Created data labels and stored data into 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