

# • Data Science Capstone Project

---

José Reimondez

<https://github.com/josereimondez29>

14/11/2024



# Outline

---



- Executive Summary (3)
- Introduction (4)
- Methodology (6)
- Results (16)
- Conclusion (46)
- Appendix (47)

# 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

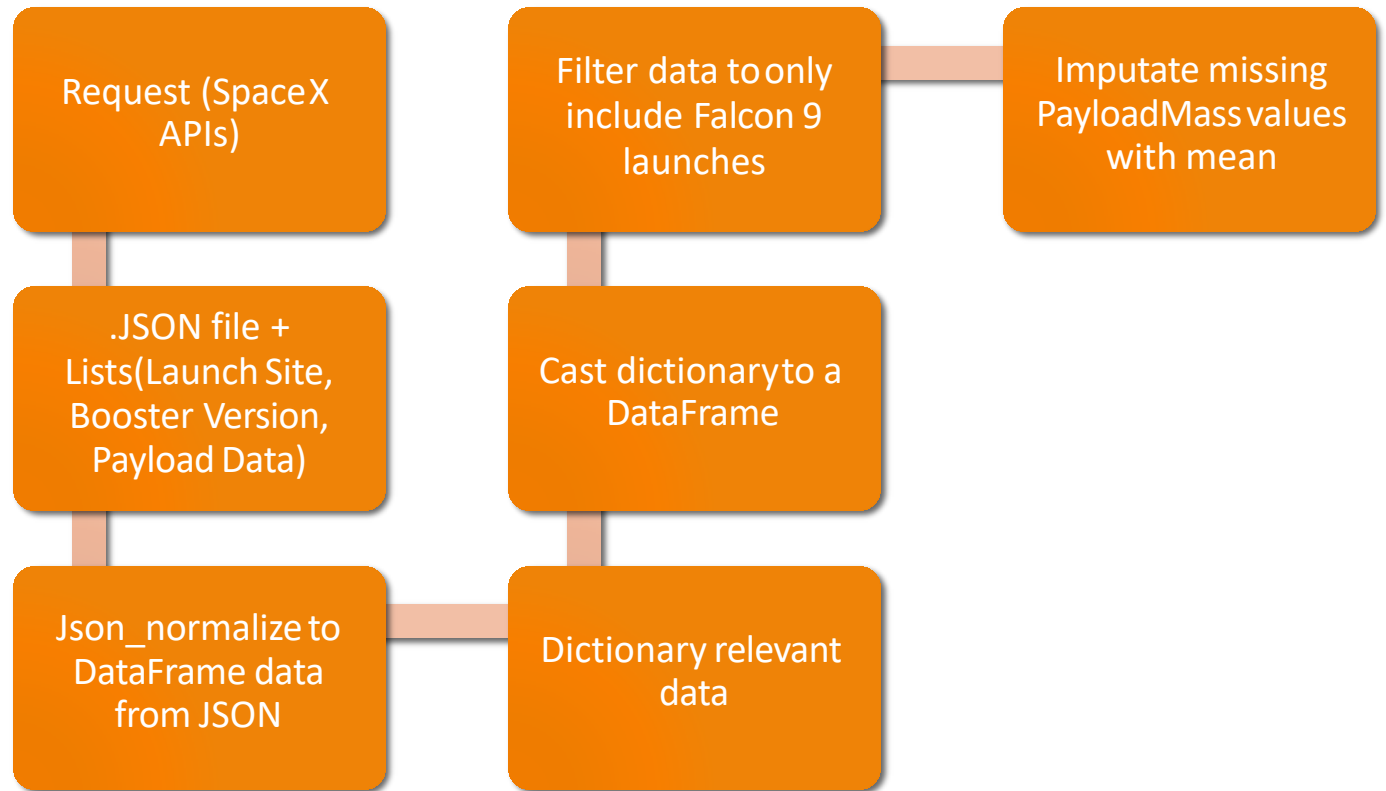
## Wikipedia Webscrape Data Columns:

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

# Data Collection— SpaceX API

GitHub url:

<https://github.com/josereimondez29/FinalBMDSC/blob/main/Data%20Collection%20Api%20.ipynb>

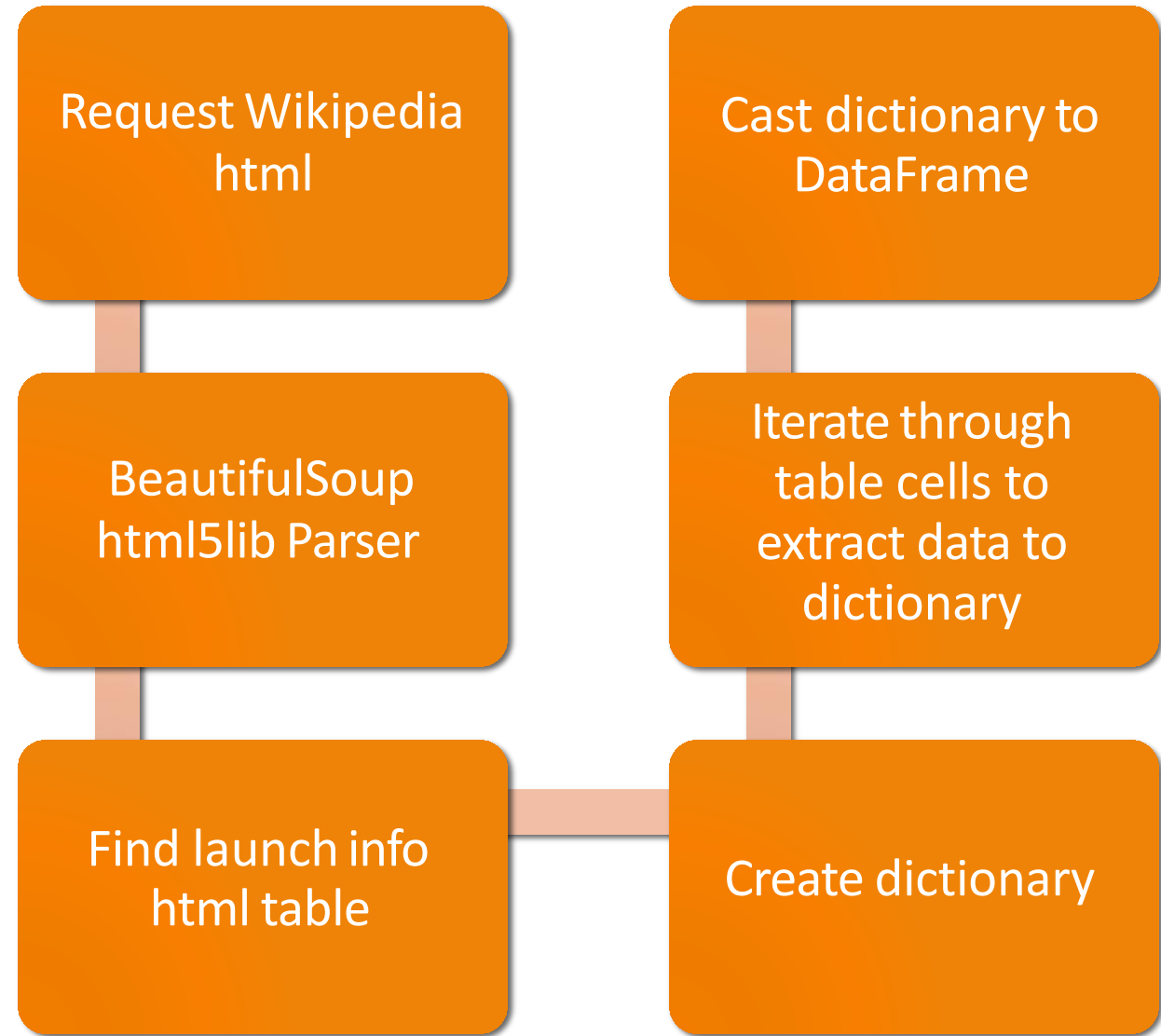




# Data Collection— Web Scraping

GitHub url:

[https://github.com/josereimondez29/  
FinalIBMDSC/blob/main/Data%20Colle  
ction%20with%20Web%20Scraping.ip  
ynb](https://github.com/josereimondez29/FinalIBMDSC/blob/main/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/josereimondez29/FinalIBMDSC/blob/main/Data%20wrangling%20.ip](https://github.com/josereimondez29/FinalIBMDSC/blob/main/Data%20wrangling%20.ipynb)

ynb

# EDA with Data Visualization

---

Exploratory Data Analysis (EDA) was performed on key variables such as Flight Number, Payload Mass, Launch Site, Orbit, Class, and Year to uncover patterns and insights. Various plots were used to visualize the relationships between these variables, including scatter plots for Flight Number vs. Payload Mass, bar plots for Flight Number vs. Launch Site, and line charts to track the Success Yearly Trend.

The goal was to identify any correlations or trends that could inform the model training process, ensuring that the chosen features have predictive value. Additionally, the analysis of the Success Rate by Orbit and the relationship between Payload and Orbit provided further insights into how different factors impact launch success, guiding feature selection for machine learning modeling. This process helped in determining which features were most influential in predicting launch outcomes, allowing for more accurate and efficient model development. By visualizing these relationships, we were able to identify outliers, trends, and data anomalies that could affect model performance.

Ultimately, EDA played a crucial role in shaping the data preparation and feature engineering steps for building a robust classification model.

# EDA with Data Visualization

Further, EDA helped uncover significant trends in the launch success rate over time, revealing that certain years exhibited higher success rates, which might be linked to technological advancements or improvements in launch practices. We also explored how launch sites affected the success of landings, finding that some sites showed higher success rates than others, potentially due to geographical or infrastructure advantages. Payload mass was found to have an interesting relationship with launch success, with certain payload categories showing distinct patterns in launch outcomes. These insights helped refine the features for the model, as we identified which variables were most predictive of success.

Additionally, the Orbit variable provided a deeper understanding of how the type of orbit influenced mission outcomes. For instance, orbits like LEO (Low Earth Orbit) and GTO (Geostationary Transfer Orbit) had different success rates, indicating the complexity of missions associated with these orbits. Identifying and visualizing these trends also led to the detection of outliers, such as extreme payload masses or unusual mission types, which were carefully handled in data preprocessing.

EDA also revealed missing data and inconsistencies that needed to be addressed, such as empty or incorrect values in some columns. These issues were resolved by either imputing missing values or dropping rows with critical missing data. Overall, EDA provided critical insights into the relationships between variables, enabling a more informed approach to feature engineering and model selection. The visualizations and statistical summaries created during this phase laid the groundwork for building a predictive model with higher accuracy and reliability.

GitHub url: <https://github.com/josereimondez29/FinalBMDSC/blob/main/EDA%20with%20Visualization.ipynb>

# 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/josereimondez29/FinalIBMDSC/blob/main/EDA%20with%20SQL.ipynb>

# Build an interactive map with Folium

---

Folium maps were used to mark key elements such as Launch Sites, successful and unsuccessful landings, and their proximity to important locations like Railway, Highway, Coast, and City. This geospatial visualization helps in understanding the strategic placement of launch sites, considering factors like accessibility, safety, and environmental constraints. By overlaying successful and unsuccessful landing data, the map provides insights into how the proximity to these key locations might influence the success of landings. This allows us to better understand the rationale behind the selection of launch site locations and their relationship to surrounding infrastructure, as well as the spatial patterns of launch success and failure.

Additionally, the proximity to critical infrastructure like highways and railways can offer operational advantages, such as easier transportation of rockets and equipment. Visualizing these relationships on a map enhances our ability to identify trends, like whether certain launch sites are more likely to experience success based on nearby factors.

This geospatial analysis provides valuable context for decision-making, helping SpaceX optimize launch site selection for future missions and predict potential challenges in different locations. By understanding how the **location of the launch site** and surrounding infrastructure affect mission success, SpaceX can improve operational efficiency and minimize logistical risks.

# Build an interactive map with Folium

Moreover, the **coastal proximity** of many launch sites is crucial for safety, as it reduces the risk to populated areas in case of a failure. Analyzing these relationships also allows for a better understanding of how **weather patterns**, **geographical features**, and **environmental regulations** impact the success of landings and launches. This data-driven approach to **location optimization** can provide strategic insights for improving the planning of future launches.

Ultimately, this analysis enhances SpaceX's ability to make more informed decisions regarding site selection, risk management, and the overall success of their missions.

GitHub url:

<https://github.com/josereimondez29/FinalBMDSC/blob/main/INTERA~1.IPY>

# 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/josereimondez29/FinalBMDSC/blob/main/spacex\\_dash\\_app.py](https://github.com/josereimondez29/FinalBMDSC/blob/main/spacex_dash_app.py)

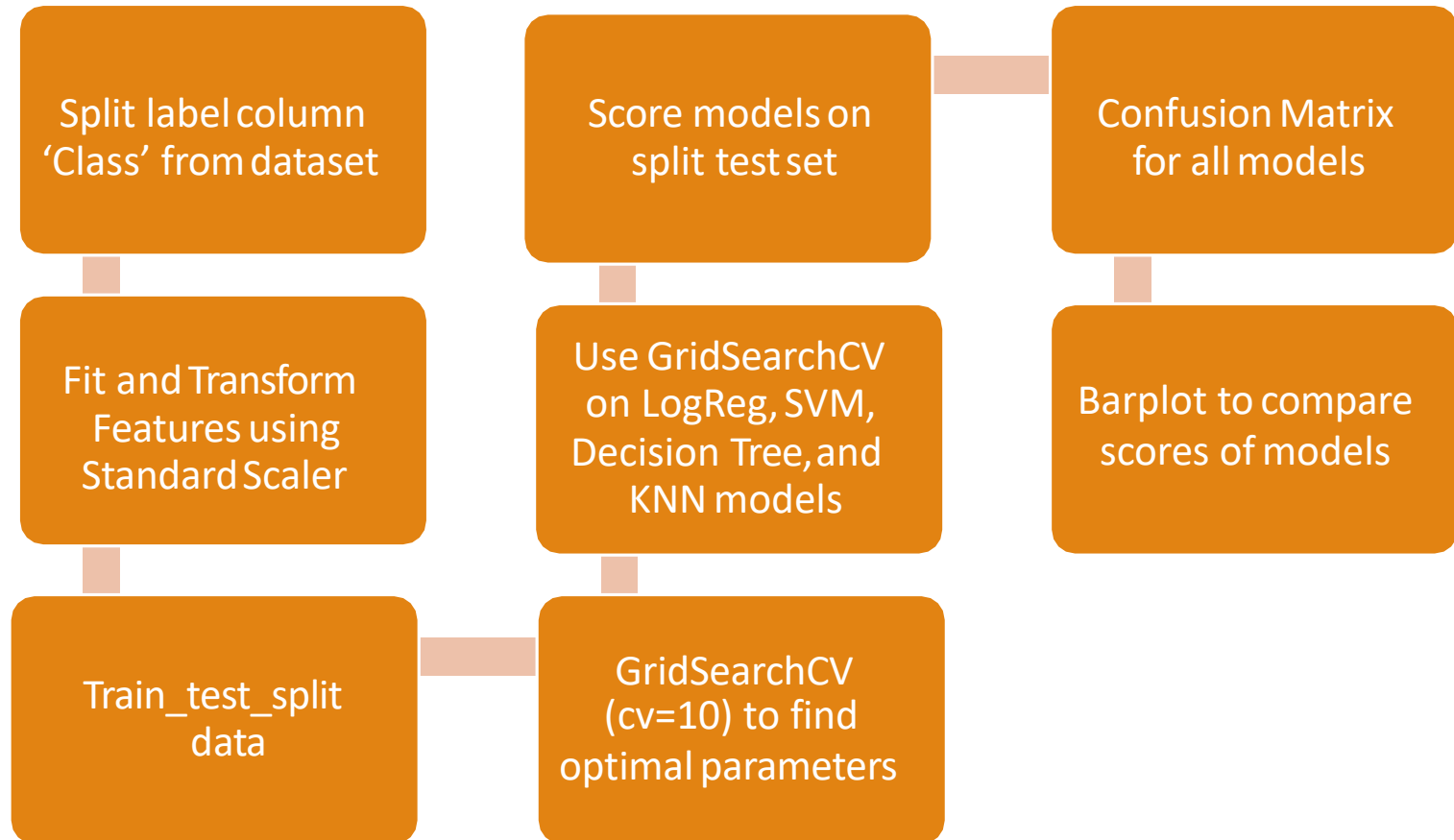


# Predictive analysis (Classification)

---

GitHub url:

<https://github.com/josereimondez29/FinalBMDSC/blob/main/Machine%20Learning%20Prediction.ipynb>

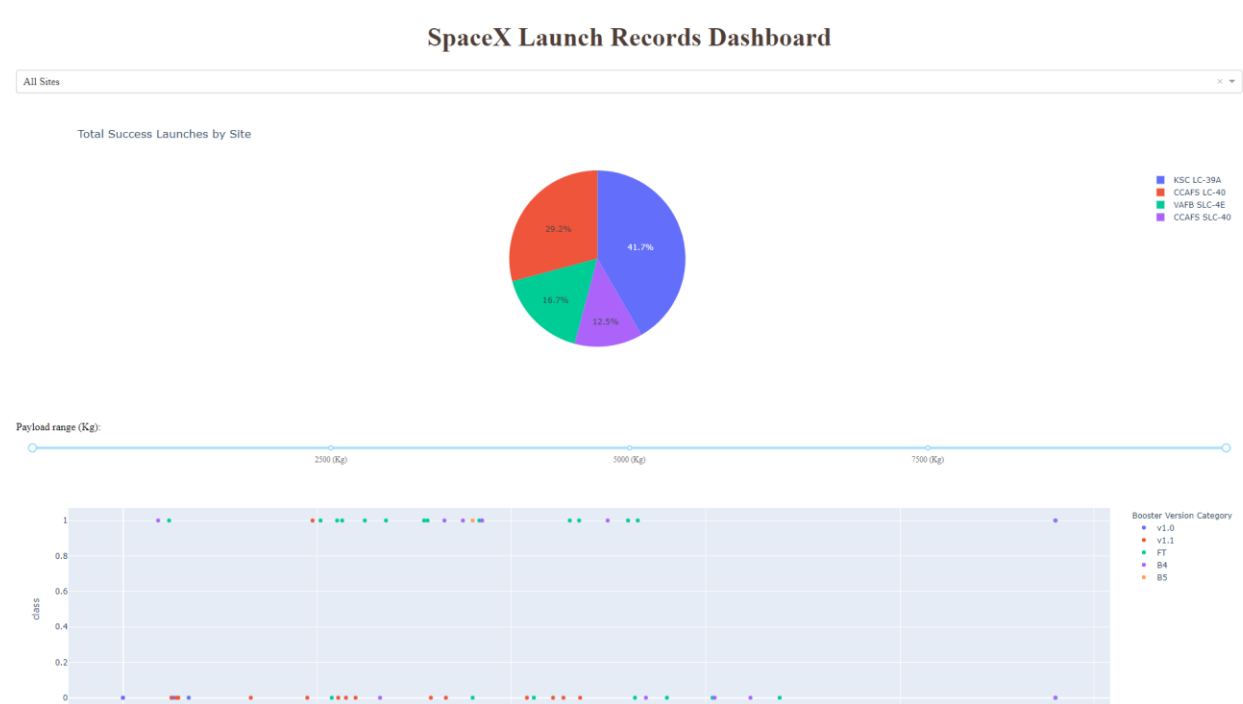


# Predictive analysis (Classification)

was performed to predict the success or failure of SpaceX landings based on various features such as **Payload Mass**, **Launch Site**, **Orbit**, and **Flight Number**. The goal was to build a model that could classify launches as successful or unsuccessful using machine learning techniques. After collecting data from the public SpaceX API and SpaceX Wikipedia page, a '**class**' label column was created to categorize the launches based on their success. The data was explored using **SQL queries**, visualizations, **Folium maps**, and interactive dashboards to identify patterns and relationships between the features and the target variable.

Categorical variables were converted to binary using **one-hot encoding**, and the data was **standardized** to ensure all features were on the same scale. **GridSearchCV** was used to fine-tune the hyperparameters for the machine learning models, optimizing their performance. Four different models were trained: **Logistic Regression**, **Support Vector Machine**, **Decision Tree Classifier**, and **K-Nearest Neighbors**.

# Results

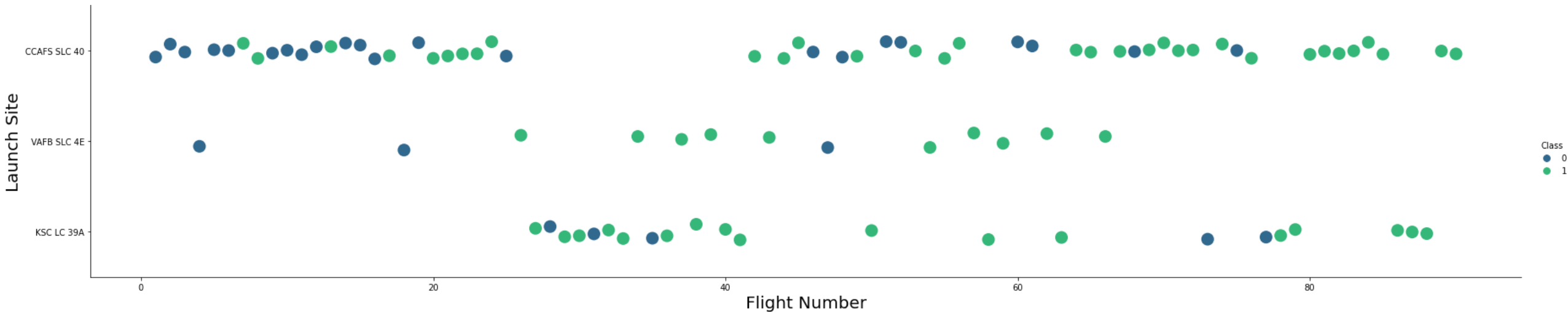


This is a preview of the Plotly dashboard. The following slides will show the results of EDA with visualization, EDA with SQL, Interactive Map with Folium, and finally the results of our model with about 83% accuracy.

# E D A with Visualization

EXPLORATORY DATA ANALYSIS WITH SEABORN PLOTS

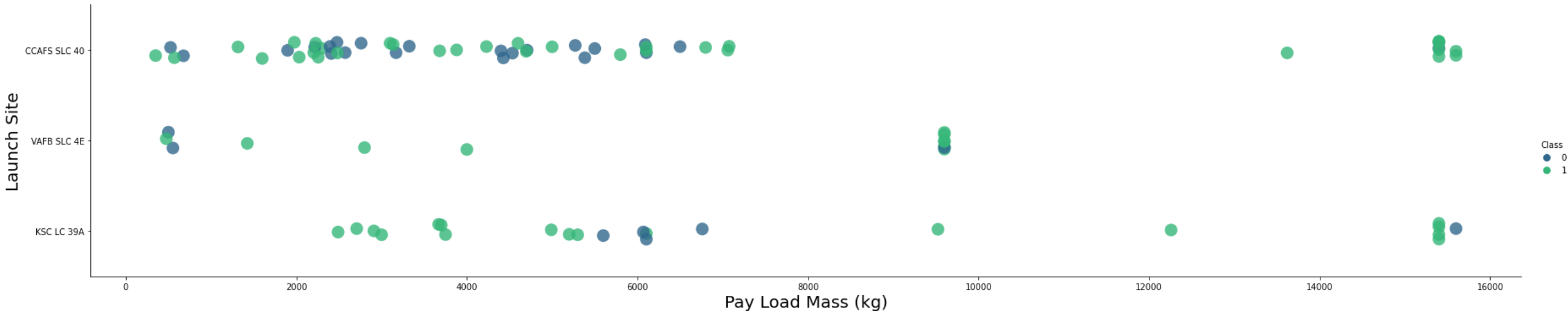
# Flight Number vs. Launch Site



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.

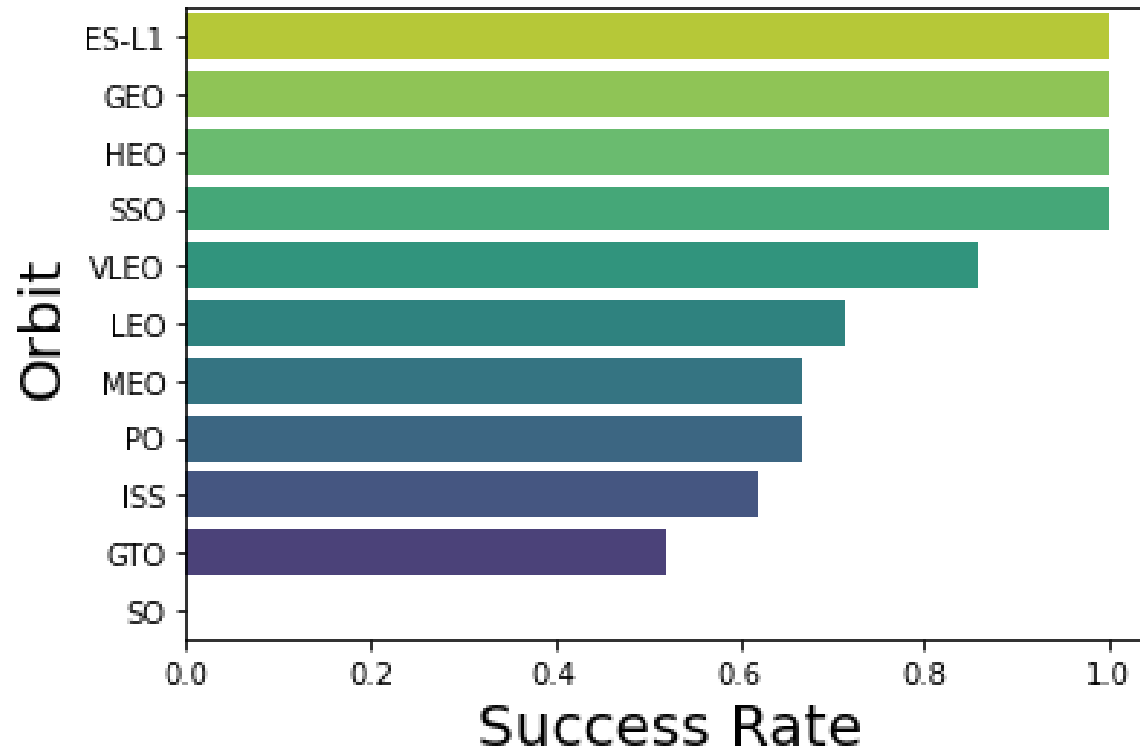
# Payload vs. Launch Site



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.

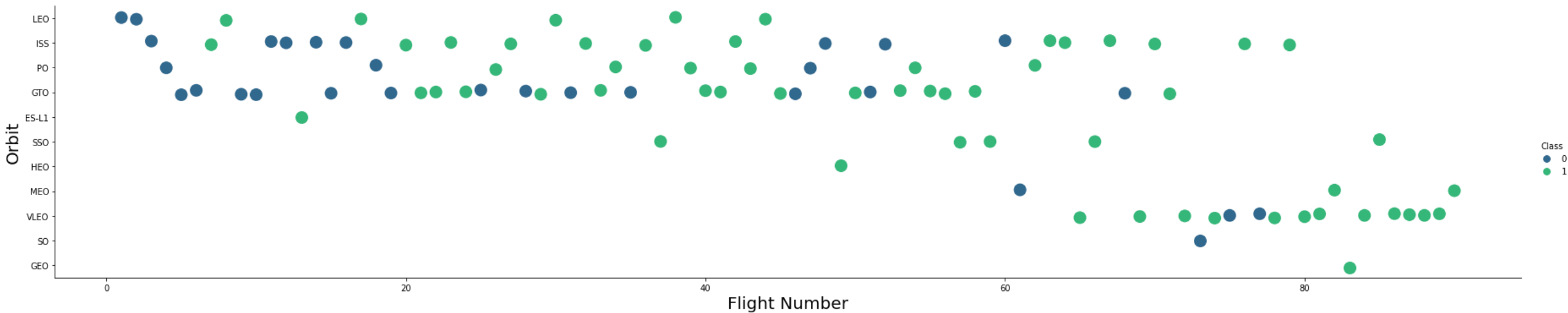
# Success rate vs. Orbit type



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



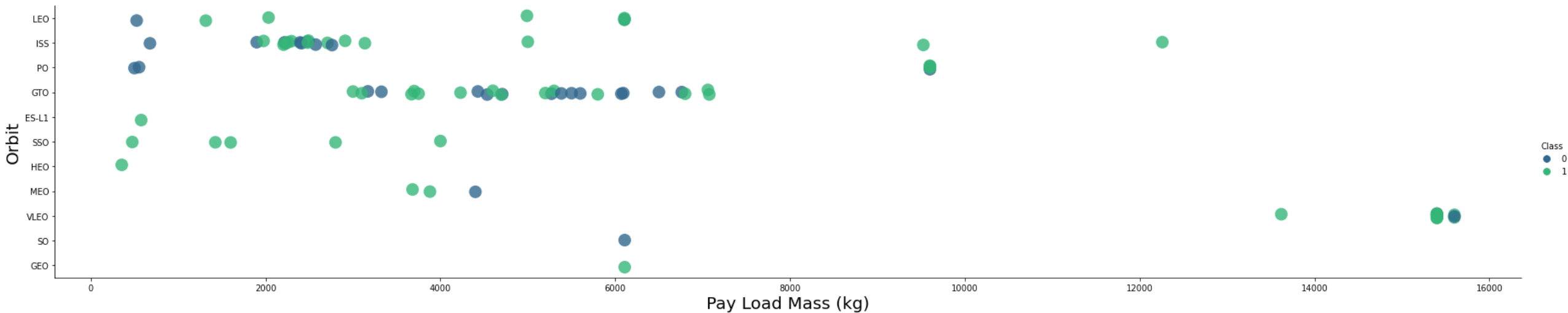
Green indicates successful launch; Purple indicates unsuccessful launch.

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



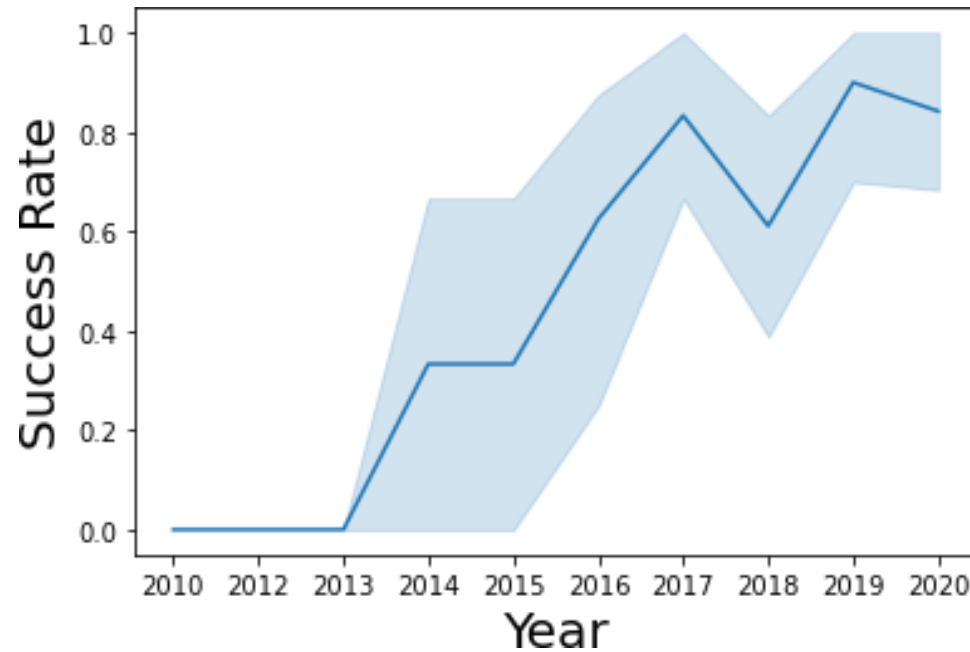
Green indicates successful launch; Purple indicates unsuccessful launch.

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%

# EDA with SQL

EXPLORATORY DATA ANALYSIS WITH SQL DB2  
INTEGRATED IN PYTHON WITH SQLALCHEMY

# 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 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 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:***@0c77d6f2-5da9-48a9-81f8-86
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 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.
```

<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.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 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:***@0c77d6f2-5da9-48a9-1
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 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 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.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.

# 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

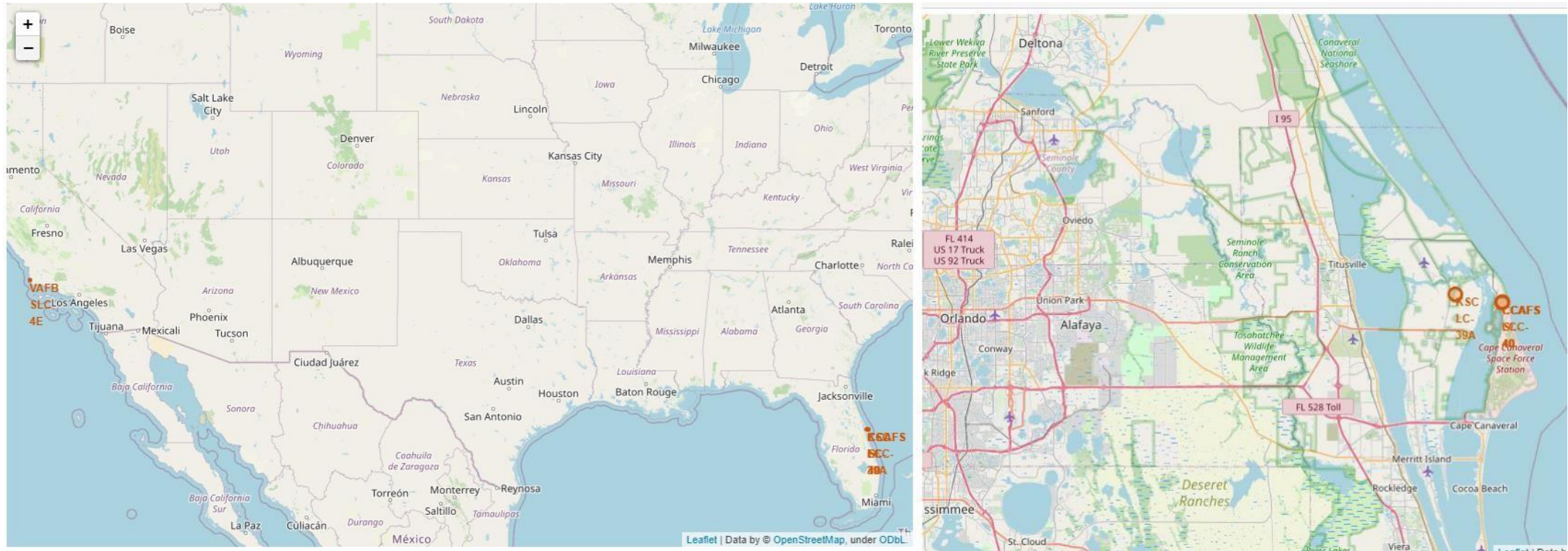
This query returns a list of successful landings and between 2010-06-04 and 2017-03-20 inclusively.

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

# Launch Site Locations

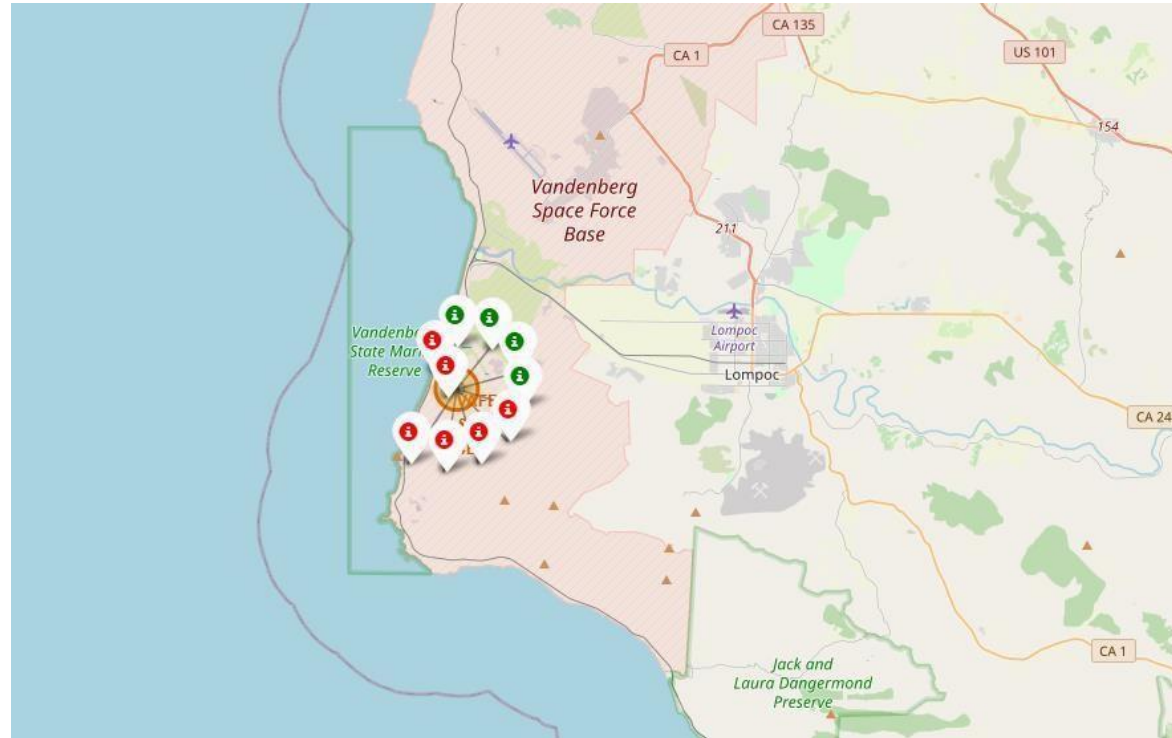


The left map shows all launch sites relative US map. The right map shows the two Florida launch sites since they are very close to each other. All launch sites are near the ocean.



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

# Build a Dashboard with Plotly Dash

# Successful Launches Across Launch Sites

---

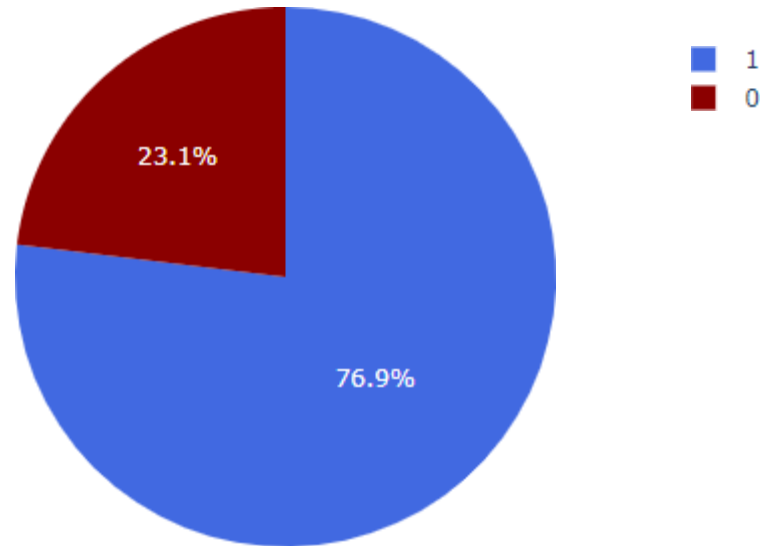


This is the distribution of successful landings across all launch sites. CCAFS LC-40 is the old name of CCAFS SLC-40 so CCAFS and KSC have the same amount of successful landings, but a majority of the successful landings were performed before the name change. VAFB has the smallest share of successful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



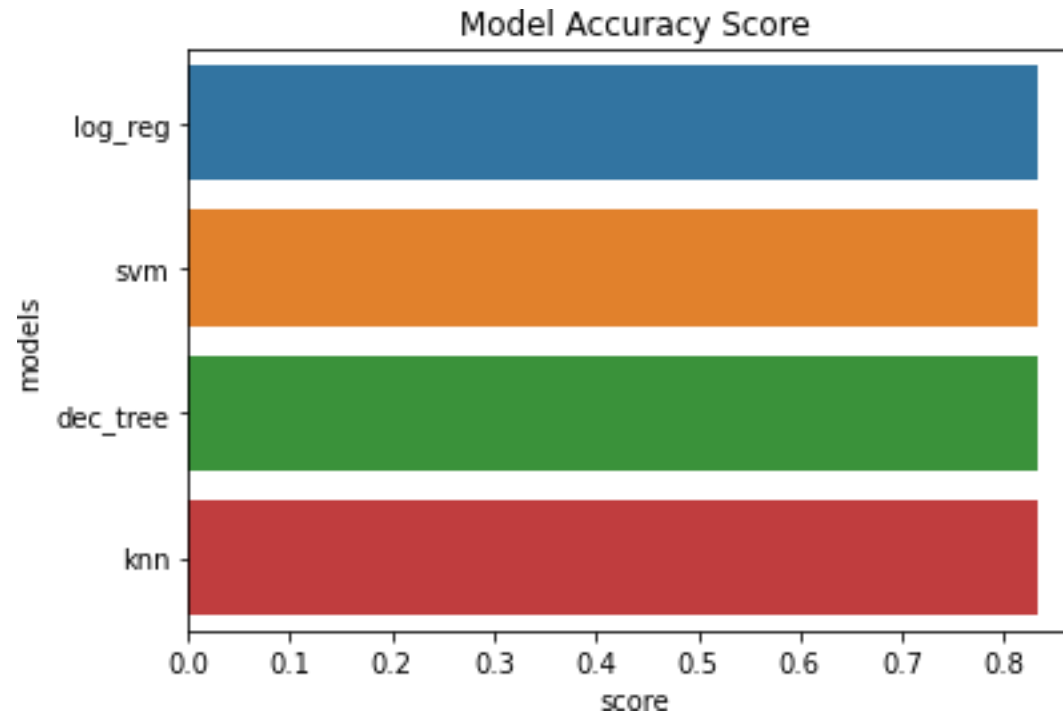
Plotly dashboard has a Payload range selector. However, this is set from 0-10000 instead of the max Payload of 15600. Class indicates 1 for successful landing and 0 for failure. Scatter plot also accounts for booster version category in color and number of launches in point size. In this particular range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

- Predictive Analysis(Classification)

---

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

# Classification Accuracy

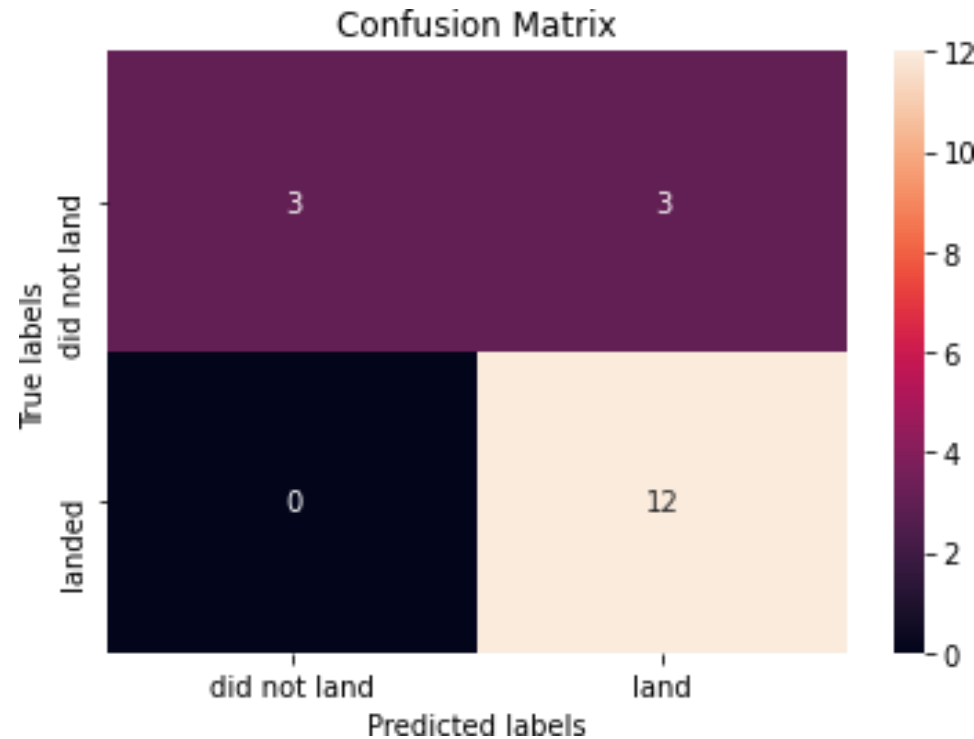


All models had virtually the same accuracy on the test set at 83.33% accuracy. It should be noted that test size is small at only sample size of 18.

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 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%
- Elon Musk 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