

# Applied Data Science Capstone Assignment

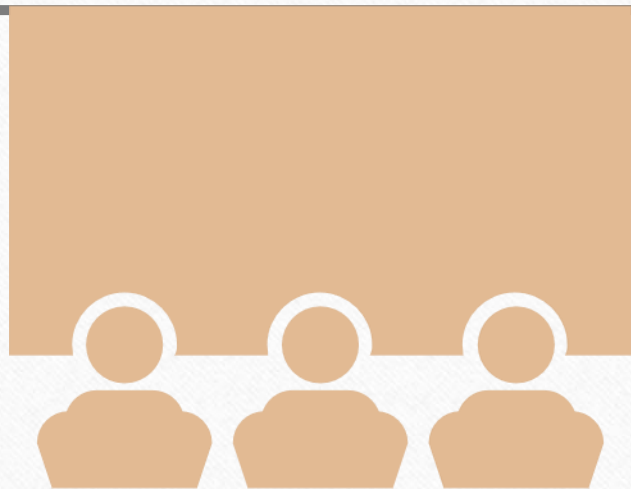
---

Satya Akhil Torlikonda

<https://github.com/SatyaAkhilTorlikonda/Course-assignment.git>

26-07-2023

# Context



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



# Summary

---

- The individual gathered data from both the public SpaceX API and SpaceX Wikipedia page. The data collected was then labeled with a column named 'class' to classify successful landings. The data exploration process involved using SQL, visualization techniques, folium maps, and dashboards to gain insights.
- Relevant columns were selected to be used as features for the machine learning models.
- The machine learning models used in the analysis were Logistic Regression, Support Vector Machine, Decision Tree Classifier, and K Nearest Neighbors. After tuning the hyperparameters using GridSearchCV, all four models performed similarly, achieving an accuracy rate of around 83.33%.

# Introduction

- 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, DASHBOARD, MODEL METHODS, VISUALIZATION AND WRANGLING.

# Data Collection Overview

---

The data collection process for this project involved a combination of two methods: API requests from SpaceX's public API and web scraping data from a table in SpaceX's Wikipedia entry.

1.API Data Collection: The first step in the data collection process was to make API requests to SpaceX's public API. This API likely provided access to a range of data related to SpaceX missions, launches, landings, and other relevant information. The data collected through these API requests was likely in a structured format, such as JSON or XML, making it easier to extract and process.

2.Web Scraping Data Collection: The second step in the data collection process was to perform web scraping on SpaceX's Wikipedia page. The goal was likely to extract specific data from a table present on the page. Web scraping involves using automated tools to navigate web pages, locate relevant data within the HTML structure, and extract the required information. In this case, the targeted table on the Wikipedia page likely contained additional information about SpaceX missions and their outcomes.



---

### Space X API Data Columns:

- FlightNumber, Date, BoosterVersion, PayloadMass, Orbit, LaunchSite, Outcome, Flights, GridFins,
- Reused, Legs, Landing Pad, Block, Reused Count, Serial, Longitude, Latitude

### Wikipedia Webscrape Data Columns:

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



# 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

# 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



# 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 payload sizes of customers and booster versions, and landing outcomes

# 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.

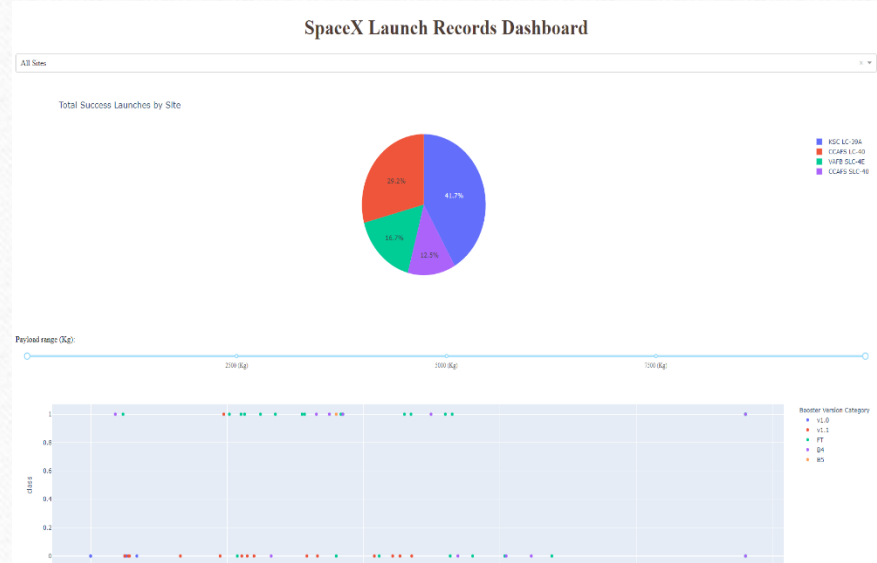


# Build a Dashboard with Plotly Dash

---

- Dashboard includes a pie chart and a scatterplot.
- 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.
- 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.

# Results



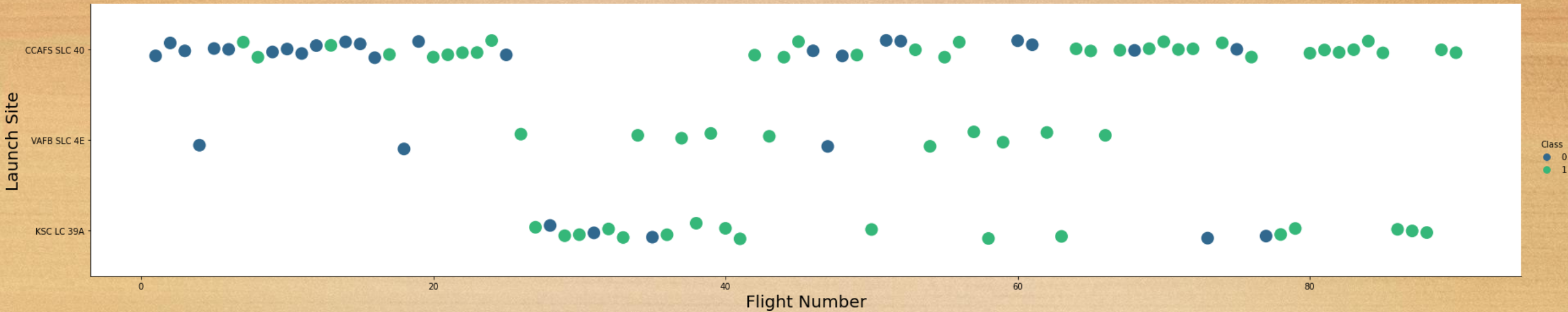
- This is a preview of the Plotly dashboard. The following slides will show the results of EDA with visualization, EDA with SQL, Interactive Map with Folium, and finally the results of our model with about 83% accuracy.



# EDA with Visualization

EXPLORATORY DATA ANALYSIS WITH SEABORN PLOTS

# Flight Number vs. LaunchSite

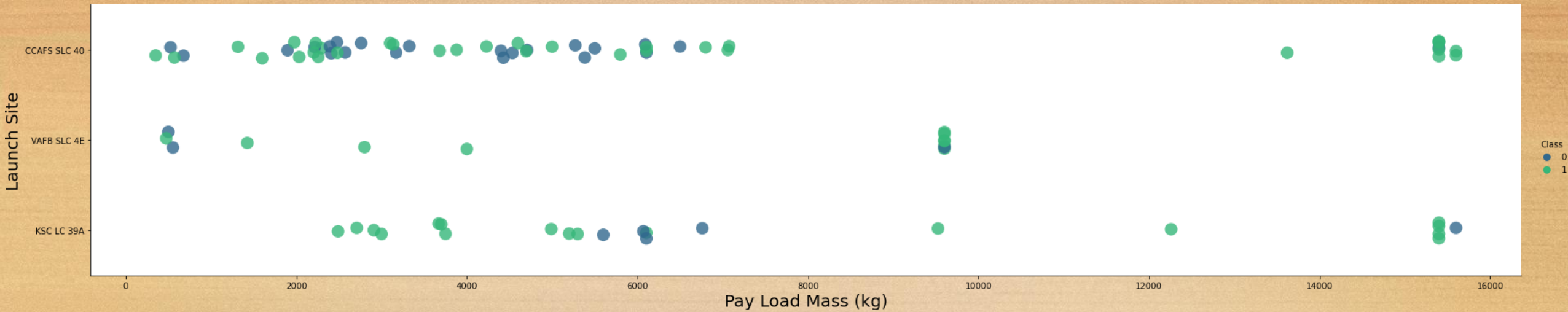


Green indicates successful launch; Purple indicates unsuccessful launch.

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.



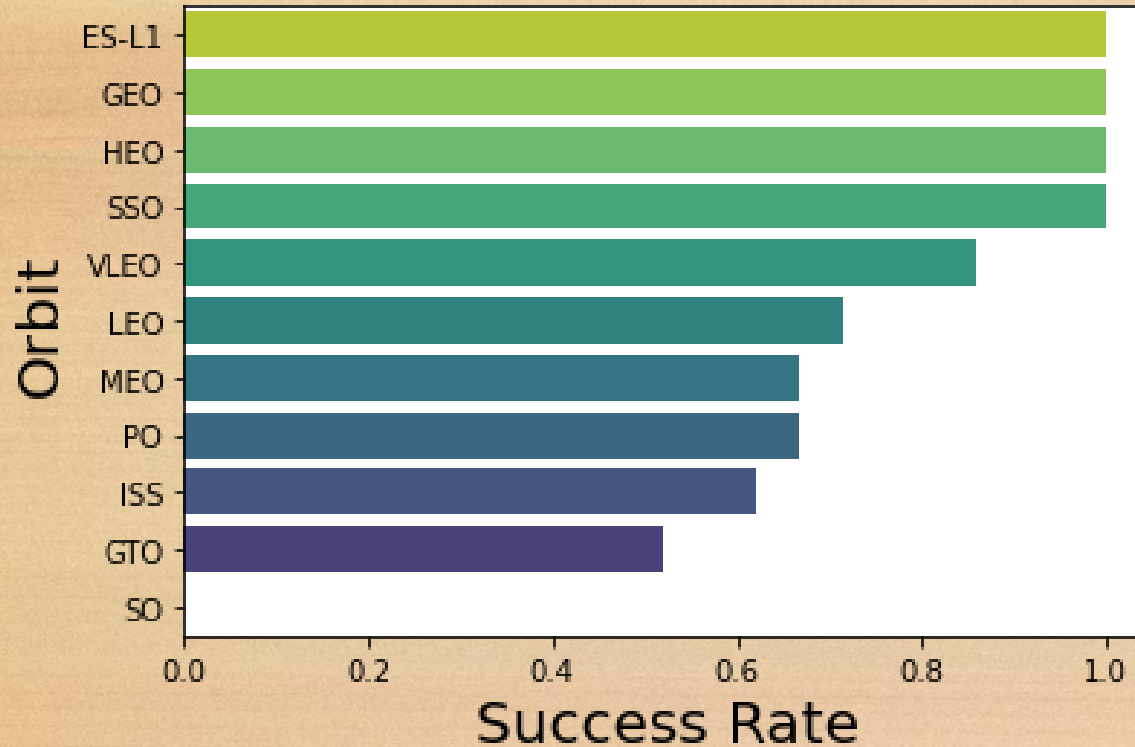
# Payload vs. Launch Site



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.

# Success rate vs. Orbit type

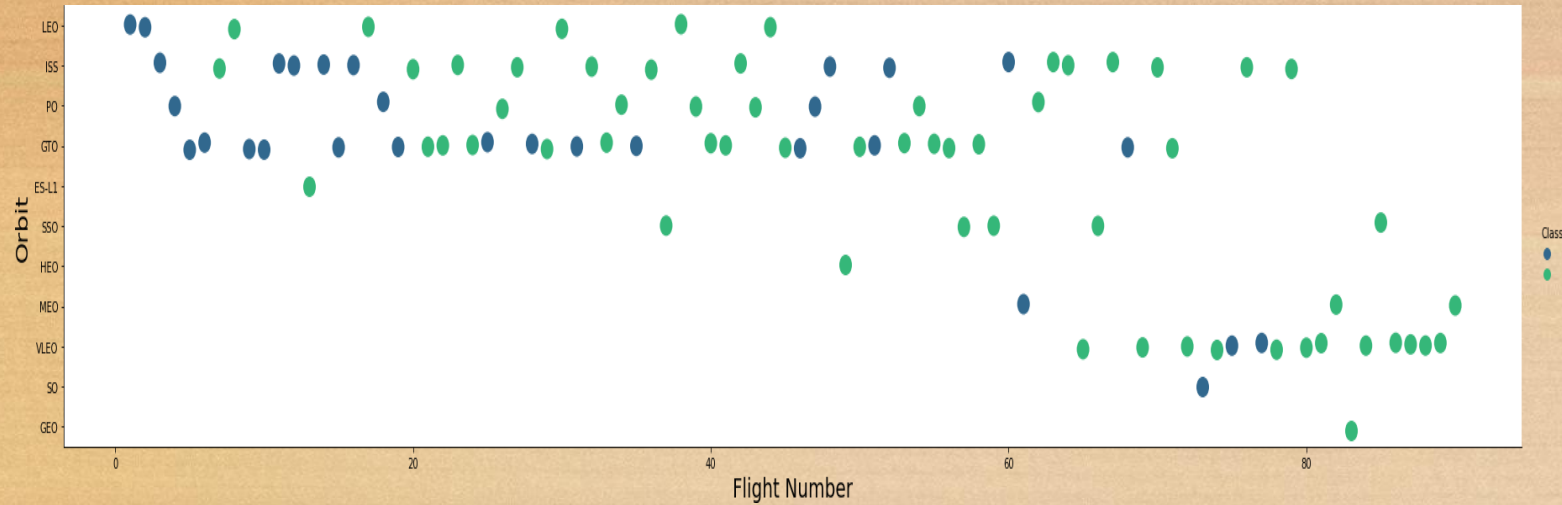


Success Rate Scale  
with 0 as 0%  
0.6 as  
60% 1 as  
100%

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. Orbit type

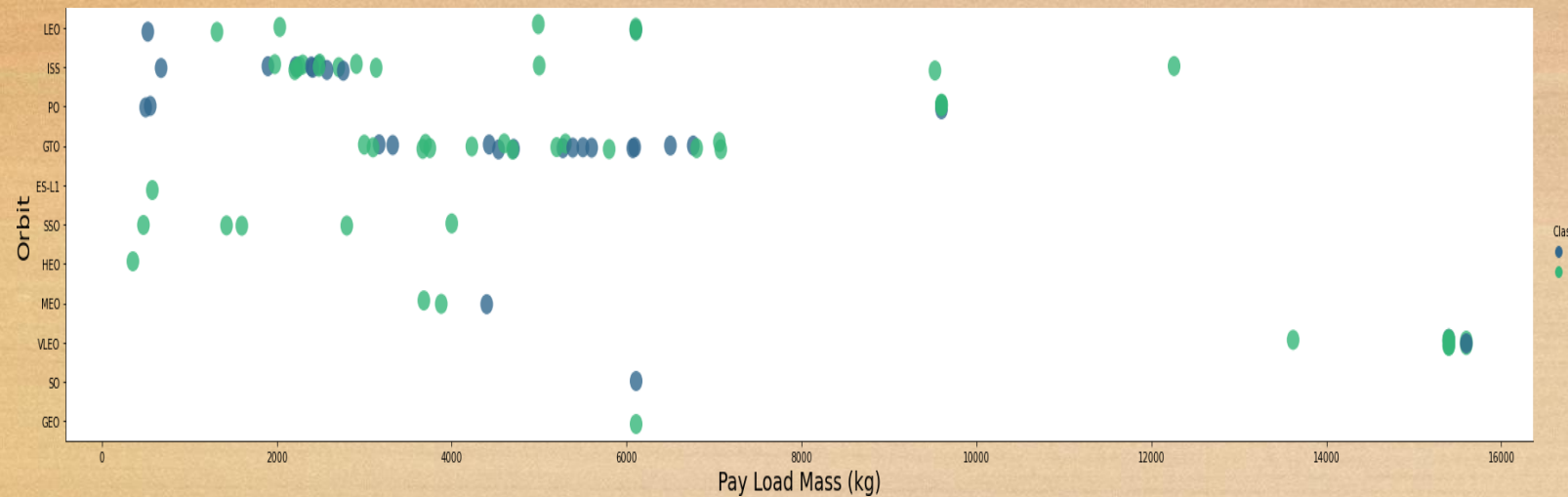


Green indicates successful launch; Purple indicates unsuccessful launch.

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.

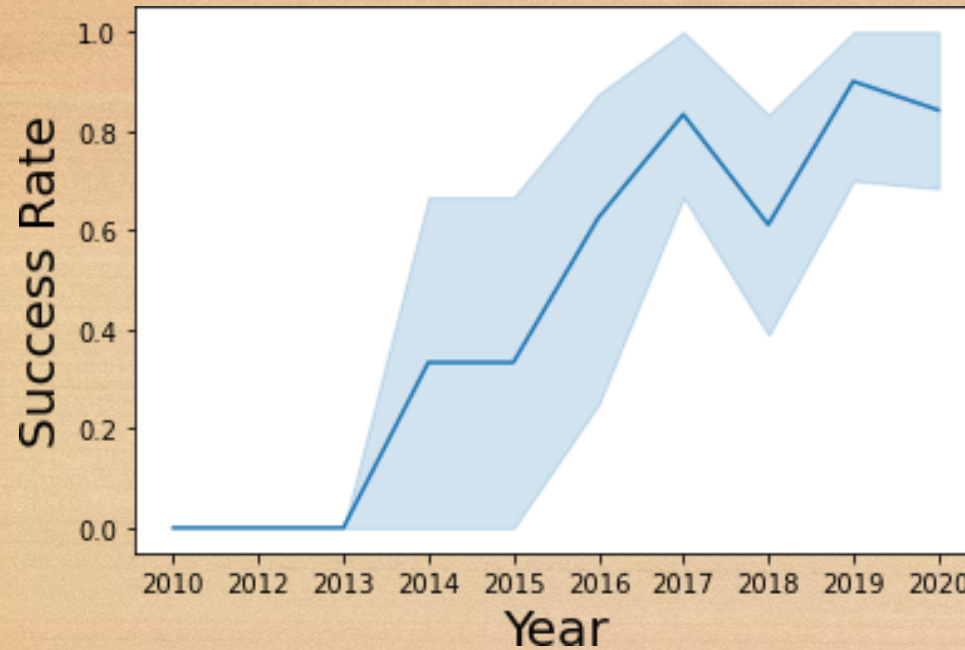
Payload mass seems to correlate with orbit

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

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

```
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  |

# Launch Site Names Beginning with 'CCA'

First five entries in database with Launch Site name 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.bs2io90108kqb1od8lcg.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          |



# Total Payload Mass from NASA

---

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).

```
%%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 |
|-------|

# Average Payload Mass by F9v1.1

---

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

```
%%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 |
|------|



# 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

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

```
%%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.database
Done.
```

| booster_version |
|-----------------|
| F9 FT B1022     |
| F9 FT B1026     |
| F9 FT B1021.2   |
| F9 FT B1031.2   |



# Total Number of Each Mission Outcome

---

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.

```
%%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          |

# Boosters that Carried Maximum Payload

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.

```
%%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.bs2io90108kqb1
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            |



# 2015 Failed Drone Ship Landing Records

---

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.

```
%%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.bs2io90108kqb1od8lcg.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 |

# Ranking Counts of Successful Landings Between 2010-06-04 and 2017-03-20

---

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

```
%%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.bs2io90108kqb1od8lç  
Done.

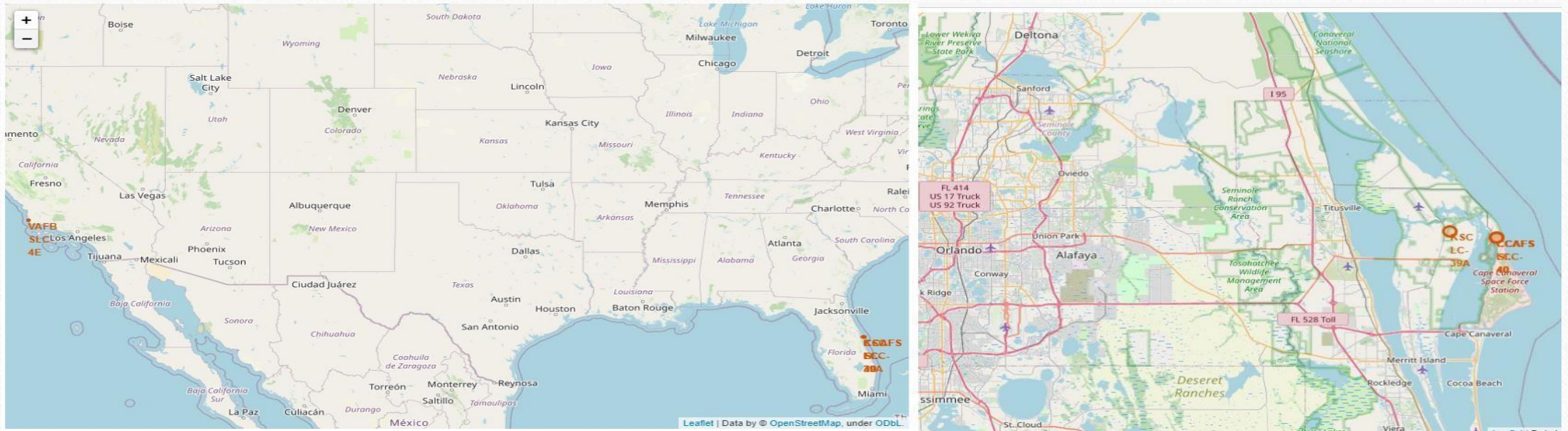
| landing_outcome      | no_outcome |
|----------------------|------------|
| Success (drone ship) | 5          |
| Success (ground pad) | 3          |



---

# Interactive Map with Folium

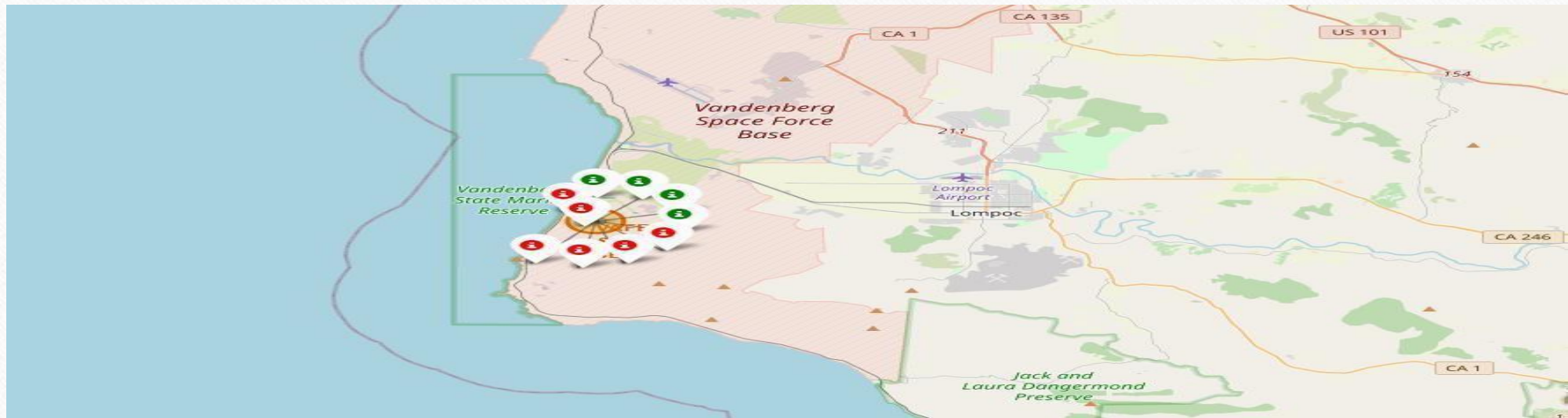
# Launch Site



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



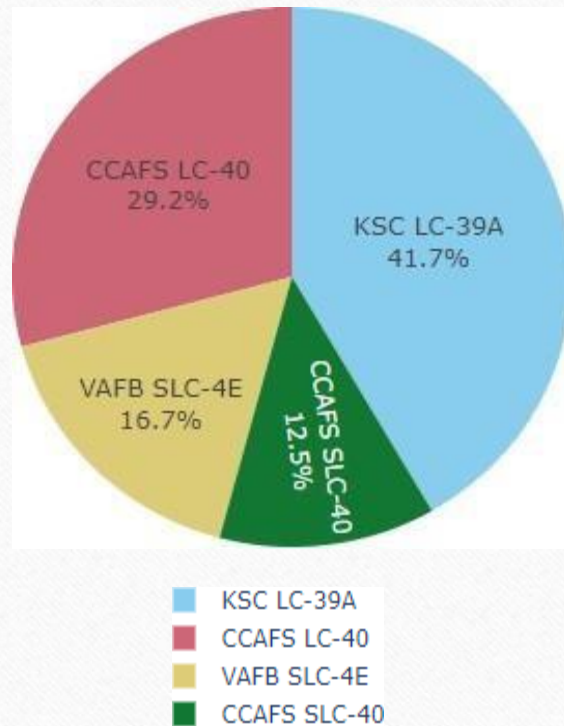
Clusters on Folium map can be clicked on to display each successful landing (green icon) and failed landing (red icon). In this example VAFB SLC-4E shows 4 successful landings and 6 failed landings.

---

# Build a Dashboard with Plotly Dash



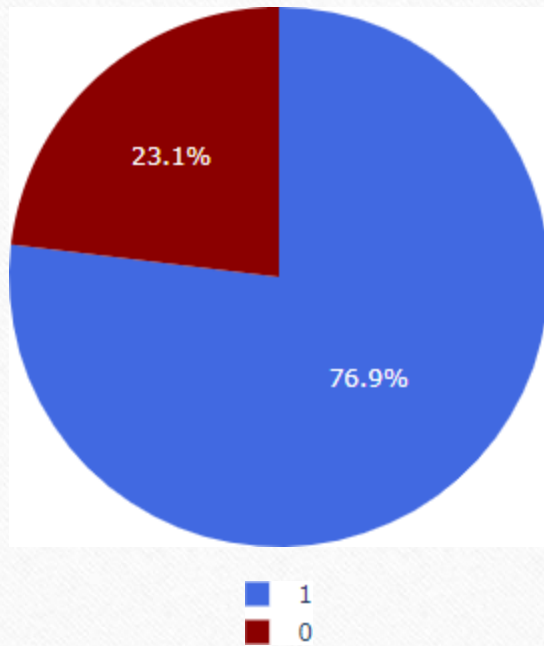
# Successful Launches Across Launch Sites



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 has the highest success rate with 10 successful landings and 3 failed landings.

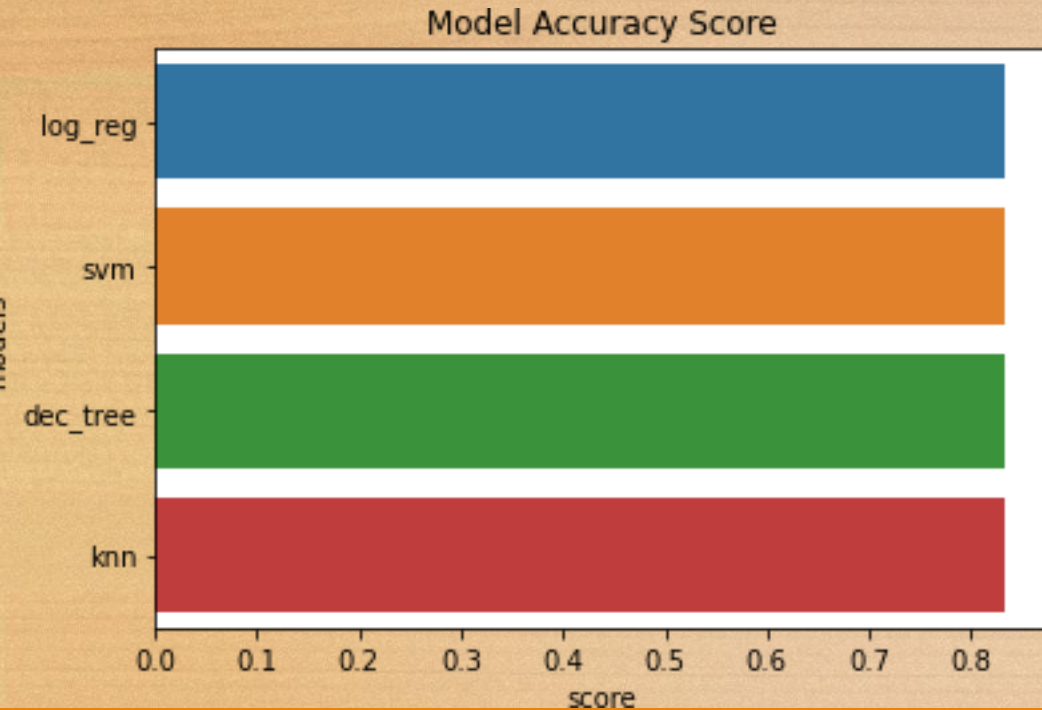


# Payload Mass vs. Success vs. Booster Version Category



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.

# Classification Accuracy



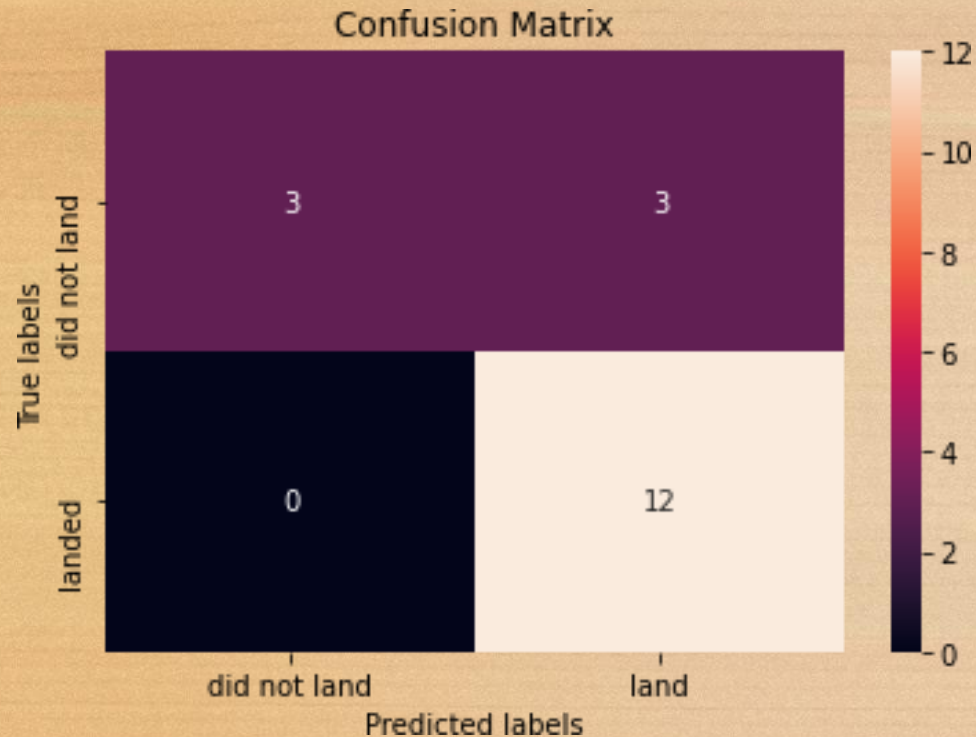
All models had virtually the same accuracy on the test set at 83.33% accuracy. It should be noted that test size is small at only sample size of 18.

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.

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

# CONCLUSION

---

- The task involved developing a machine learning model for a company named Space Y, which aims to bid against SpaceX. The specific goal of the model was to predict whether the Stage 1 of a rocket launch by SpaceX will successfully land, which could potentially save around \$100 million USD in costs.
- To build the model, data was collected from both a public SpaceX API and by web scraping information from SpaceX's Wikipedia page. The collected data was labeled appropriately and stored in a DB2 SQL database for easy access and manipulation.
- A dashboard was created to visualize the data, providing insights into the performance of SpaceX launches and landings. Using the collected data, a machine learning model was developed, which achieved an accuracy of 83% in predicting successful Stage 1 landings.



- The model's accuracy is considered relatively high, and Space Y's CEO, Alon Mask, can use this model to make more informed decisions about whether to proceed with a launch or not. By predicting the success of a Stage 1 landing, they can avoid potential losses of \$100 million USD.
  - To further improve the model's accuracy, the suggestion is to gather more data. Additional data can help in better understanding the patterns and factors that contribute to successful landings, leading to the selection of the best machine learning model for the task.
- 
- Overall, the machine learning model developed in this project provides Space Y with a valuable tool to assess the likelihood of a successful Stage 1 landing, allowing them to make more strategic and cost-effective decisions in their bid against SpaceX.

# APPENDIX

---

GitHub repository url:

<https://github.com/SatyaAkhilTorlikonda/Coursera-assignment.git>

Instructors:

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

**Grateful and Thanks to All Instructors.**