

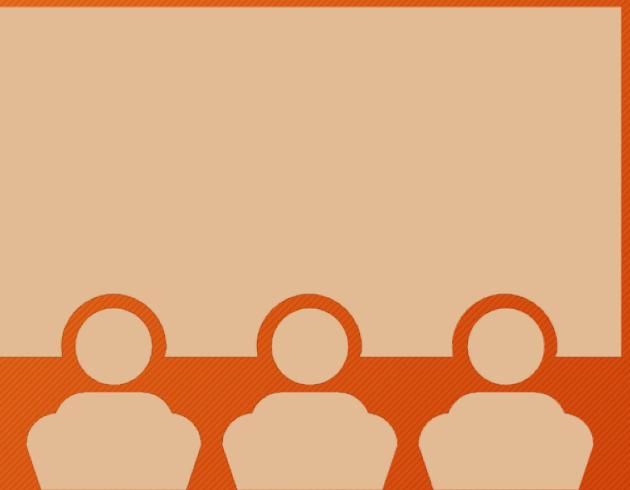
- Data Science Capstone Project - Coursera

Walid Alakk

<https://github.com/Walid-Alakk>

13/07/2024

Outline



- Executive Summary (3)
- Introduction (4)
- Methodology (6)
- Results (16)
- Conclusion (45)

Executive Summary

The commercial space age is advancing rapidly, with companies like SpaceX leading the way in affordable space travel. SpaceX's success is partly due to the reusability of their Falcon 9 rocket's first stage, which significantly reduces launch costs. In this project, you will take on the role of a data scientist working for a new rocket company, Space Y, founded by Allon Musk. Your mission is to determine the cost of each launch and predict the success of the Falcon 9's first stage landing.



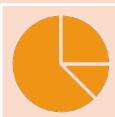
Data Collection and Preprocessing:

Gather historical data on SpaceX Falcon 9 launches.
Clean and preprocess the data for accurate analysis.



Exploratory Data Analysis (EDA):

Perform EDA to uncover patterns in launch data.
Visualize metrics like launch success rates and payload details.



Dashboard Creation:

Create interactive dashboards to visualize findings.
Include visualizations like line charts, bar charts, and pie charts.



Machine Learning Model Development:

Develop and train a model to predict first stage landing success.
Use classification algorithms to make predictions.



Model Evaluation and Deployment:

Evaluate the model using metrics like accuracy and F1-score.
Deploy the model for real-time predictions.

Introduction



PERFECTING
PROPULSIVE
LANDING

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

- 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 DATACOLLECTION, WRANGLING, VISUALIZATION,
DASHBOARD, AND MODEL METHODS

Data Collection Overview

In this capstone assignment, we focus on collecting and processing SpaceX launch data using the SpaceX REST API. This API provides detailed information about rocket launches, including rocket specifications, payload details, and landing outcomes. We retrieve past launch data from the endpoint `api.spacexdata.com/v4/launches/past` by performing a GET request with Python's `requests` library. The response is a JSON list of objects, which we convert into a flat table format using the `json_normalize` function.

In addition to the API, we use web scraping to gather Falcon 9 launch data from related Wikipedia pages. Using the `BeautifulSoup` package, we scrape HTML tables and convert them into Pandas dataframes for further analysis. This complementary method ensures that we have a comprehensive dataset for our analysis.

Data wrangling is a crucial step in preparing the dataset for analysis. We handle identification numbers by making additional API calls to retrieve detailed data. We filter out Falcon 1 launches, focusing only on Falcon 9 data. Additionally, we address missing values in the `PayloadMass` column by calculating the mean and replacing NULL values. The `LandingPad` column's NULL values will be handled using one-hot encoding later on.

The project aims to prepare a clean and comprehensive dataset for predictive analysis. By collecting, cleaning, and preprocessing the SpaceX launch data, we set the stage for training machine learning models to predict whether SpaceX will attempt to land a rocket. This work combines data collection, web scraping, and data wrangling to build a robust foundation for further analysis.

Data Collection – SpaceX API

GitHub url:

https://github.com/navassherif98/IBM_Data_Science_Professional_Certification/blob/master/10.Applied_Data_Science_Capstone/Week%201%20Introduction/Data%20Collection%20Api%20.ipynb



Data Collection – Web Scraping

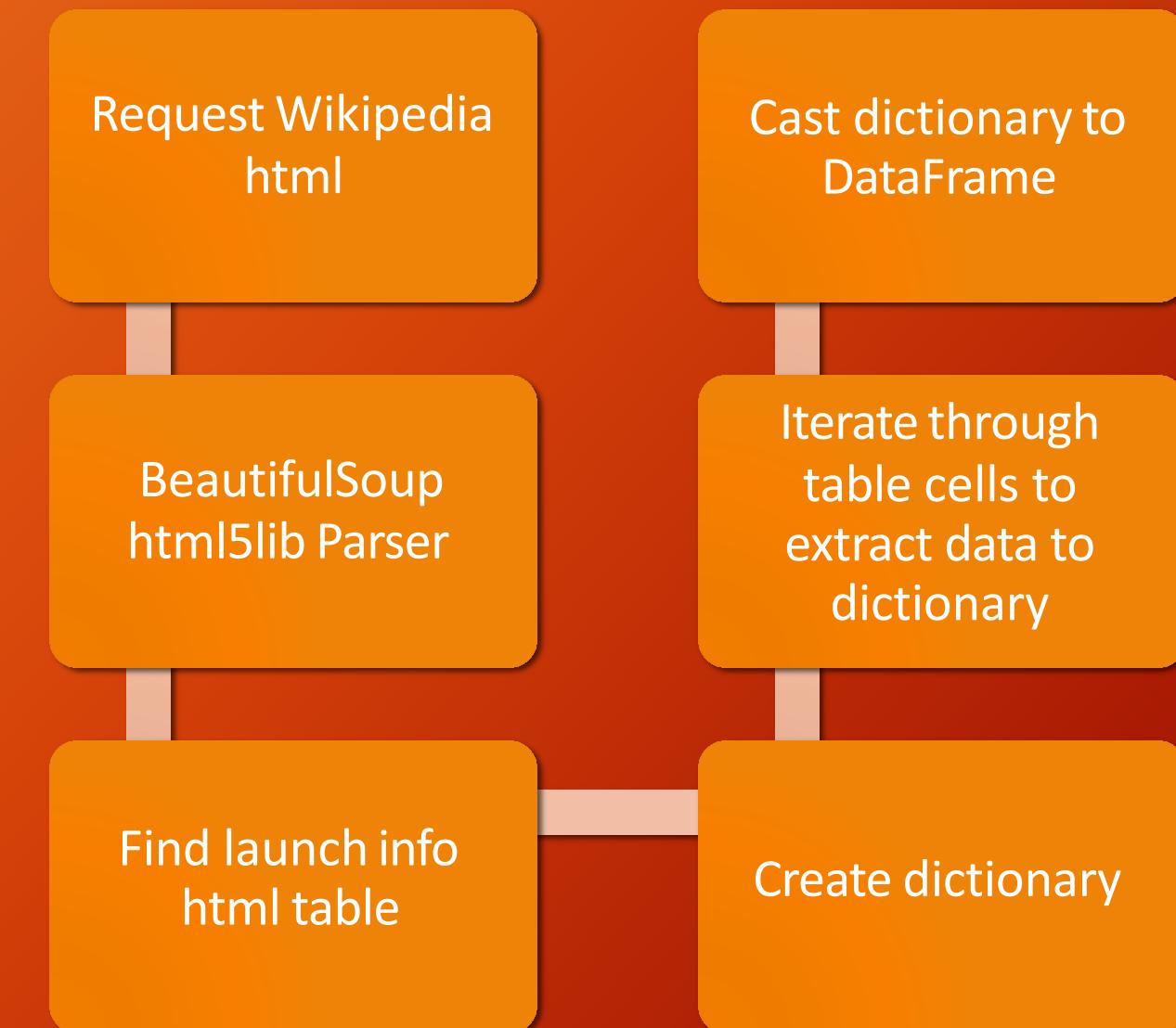
GitHub url:

Data Collection

<https://github.com/Walid-Alakk/Data-Science-Capstone-Project-SpaceX/blob/main/jupyter-labs-spacex-data-collection-api.ipynb>

Web Scraping

<https://github.com/Walid-Alakk/Data-Science-Capstone-Project-SpaceX/blob/main/jupyter-labs-webscraping.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/Walid-Alakk/Data-Science-Capstone-Project-SpaceX/blob/main/labs-jupyter-spacex-Data%20wrangling.ipynb>

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/Walid-Alakk/Data-Science-Capstone-Project-SpaceX/blob/main/EDA%20with%20Visualization%20Lab.ipynb>

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/Walid-Alakk/Data-Science-Capstone-Project-SpaceX/blob/main/jupyter-labs-eda-sql-coursera_sqlite.ipynb

Build an interactive map with Folium

13

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/Walid-Alakk/Data-Science-Capstone-Project-SpaceX/blob/main/Interactive%20Visual%20Analytics%20with%20Folium%20lab.ipynb>

Build a Dashboard with PlotlyDash

14

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:

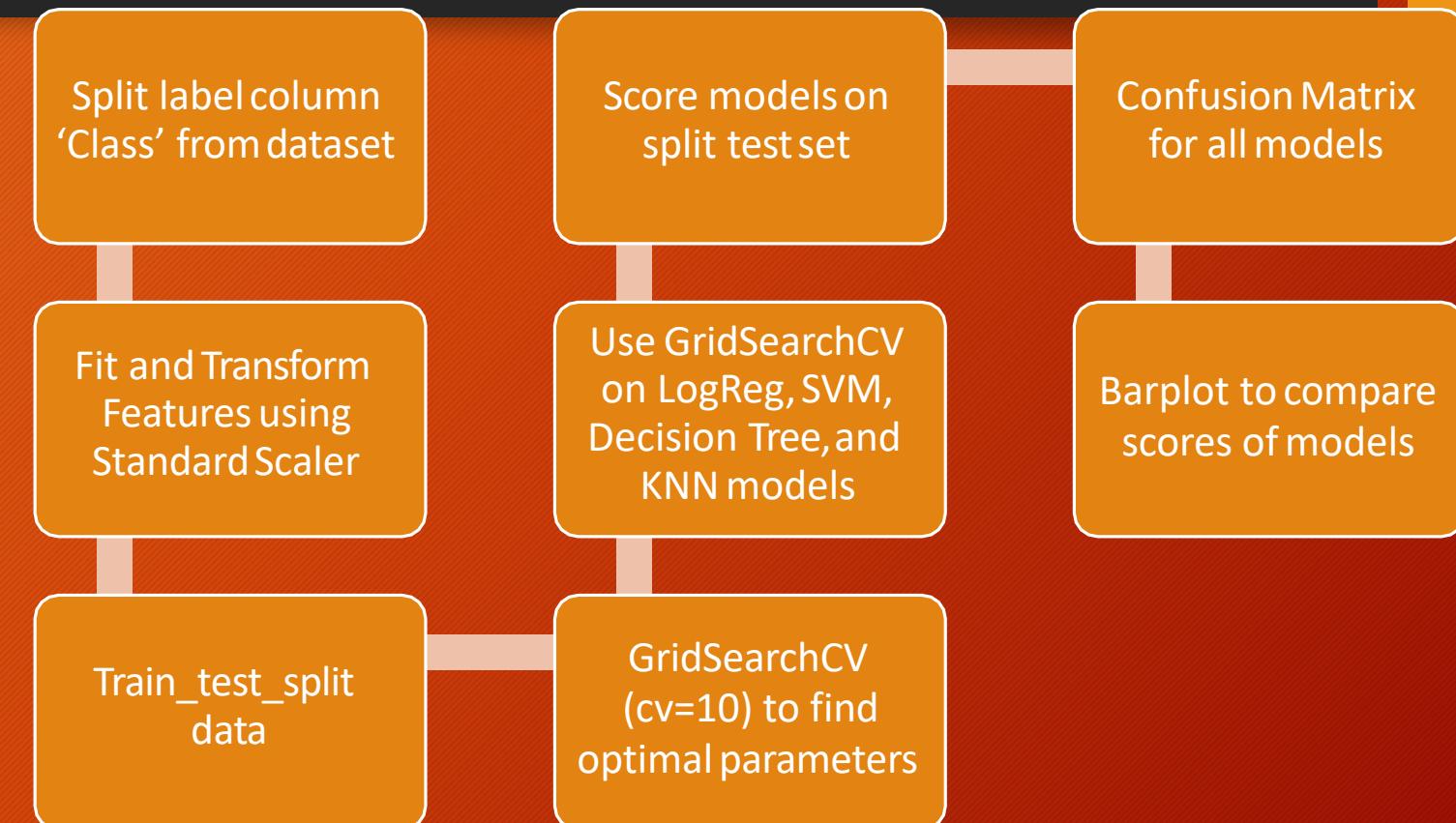
https://github.com/Walid-Alakk/Data-Science-Capstone-Project-SpaceX/blob/main/spacex_dash_app.py

Predictive analysis (Classification)

15

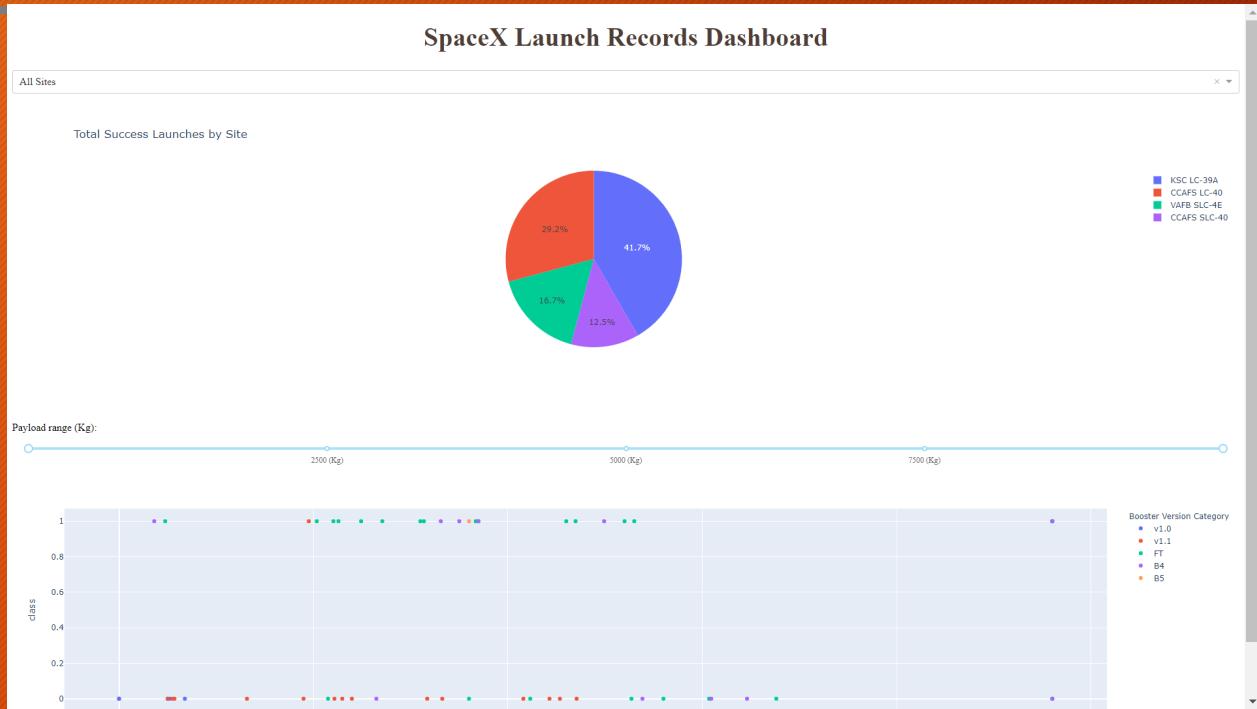
GitHub url:

https://github.com/Walid-Alakk/Data-Science-Capstone-Project-SpaceX/blob/main/SpaceX_Machine%20Learning%20Prediction_Part_5.ipynb



Results

16



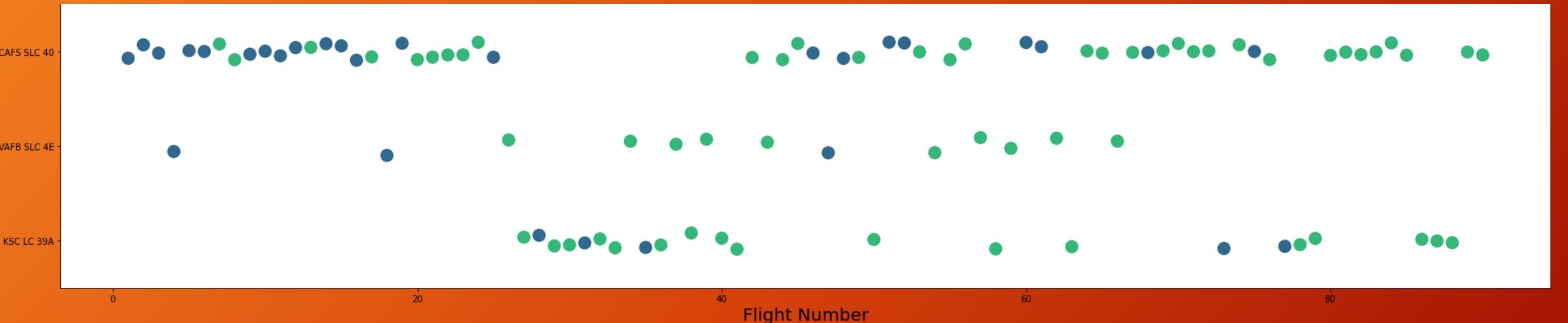
This is a preview of the Plotly dashboard. The following slides will show the results of EDA with visualization, EDA with SQL, Interactive Map with Folium, and finally the results of our model with about 83% accuracy.

EDA with Visualization

EXPLORATORY DATA ANALYSIS WITH SEABORN PLOTS

Flight Number vs. LaunchSite

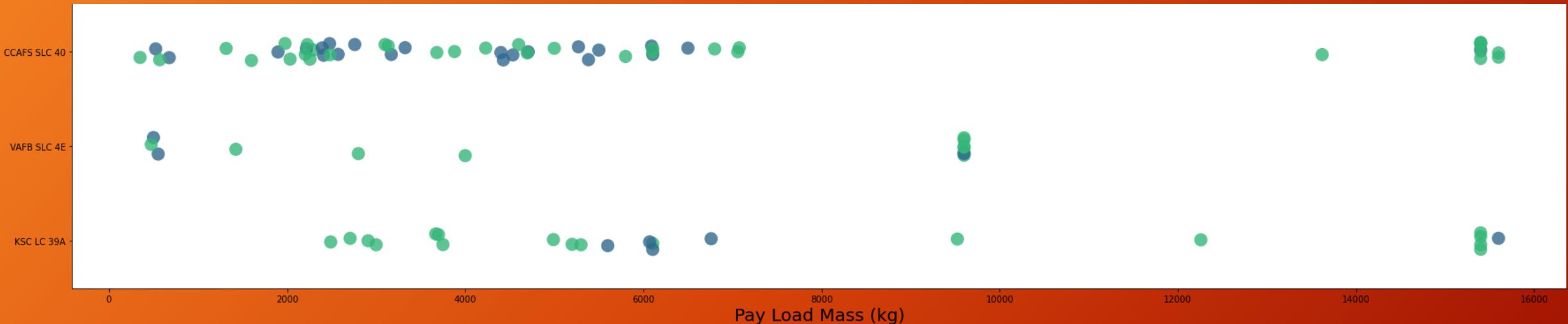
18



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

19



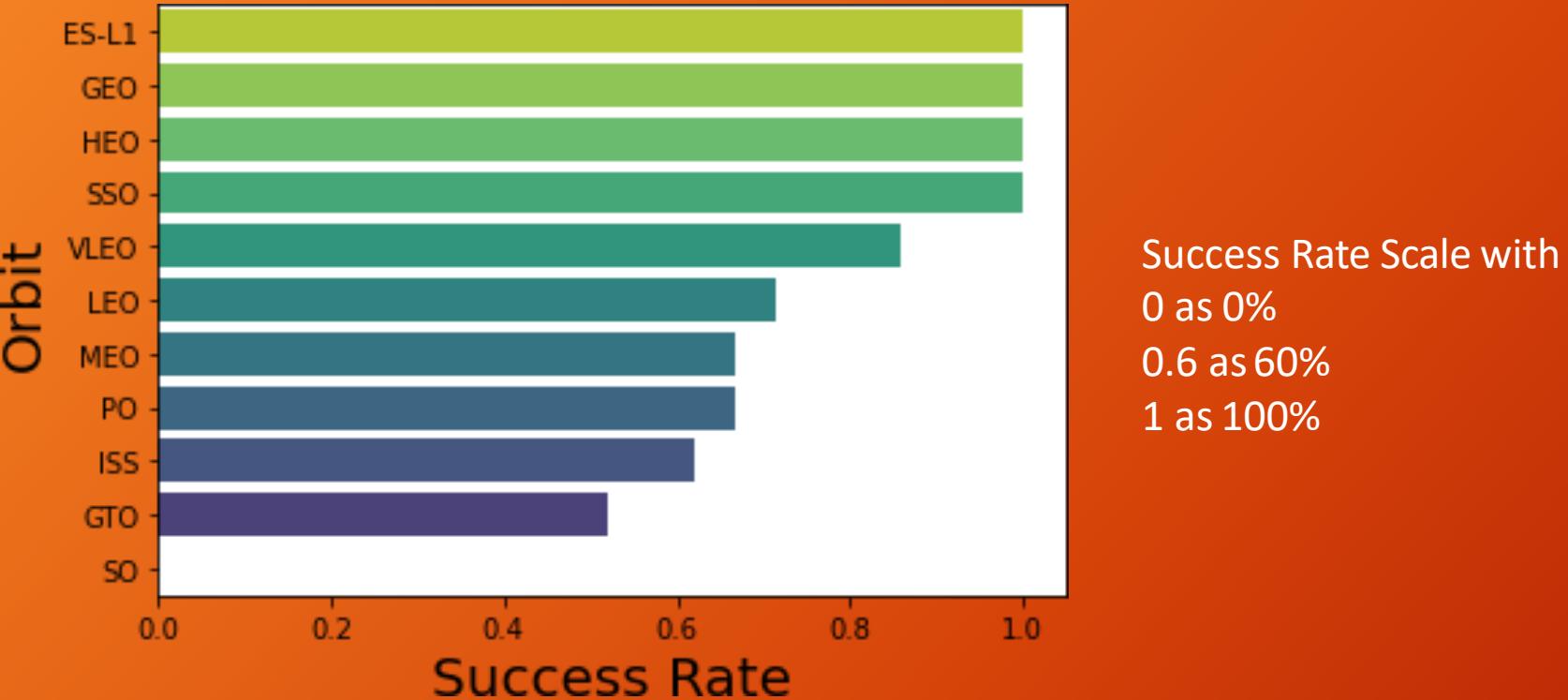
Green indicates successful launch; Purple indicates unsuccessful launch.

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

20



ES-L1 (1), GEO (1), HEO (1) have 100% success rate (sample sizes in parenthesis) SSO (5) has 100% success rate

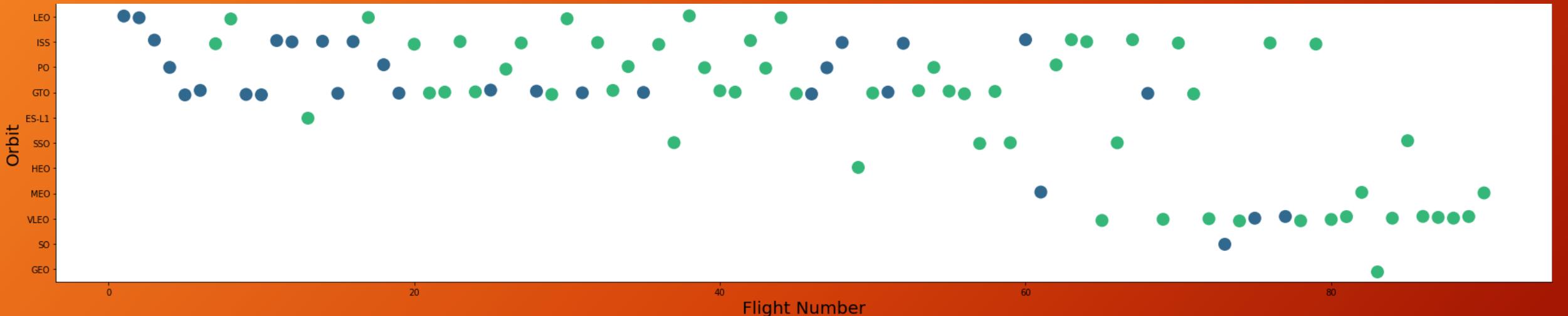
VLEO (14) has decent success rate and attempts

SO (1) has 0% success rate

GTO (27) has the around 50% success rate but largest sample

Flight Number vs. Orbit type

21



Green indicates successful launch; Purple indicates unsuccessful launch.

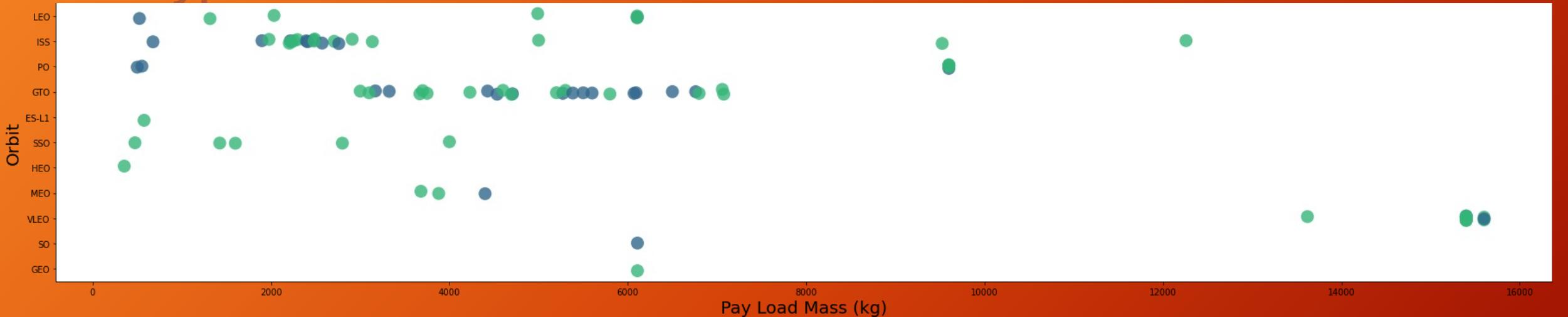
Launch Orbit preferences changed over Flight Number.

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.

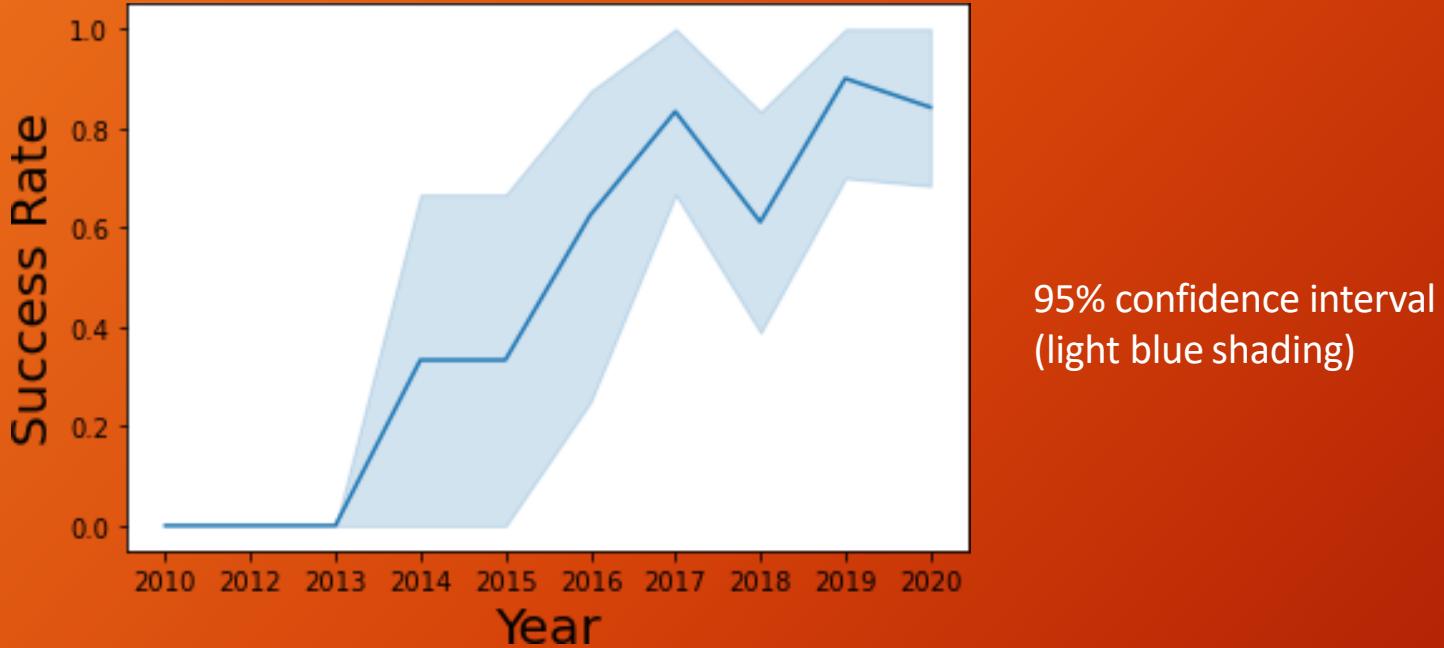
Payload mass seems to correlate with orbit

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

23



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

All Launch Site Names

25

```
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
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

First five entries in database with Launch Site name beginning with CCA.

Total Payload Mass from NASA

27

```
%%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

28

```
%%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

29

```
%%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

30

```
%%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 EachMissionOutcome

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:***@0c77d6f2-5da9-48a9-8
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.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

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.bs2io90108kqb1od81cg.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

34

```
%%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.bs2io90108kqb1od81cg
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.

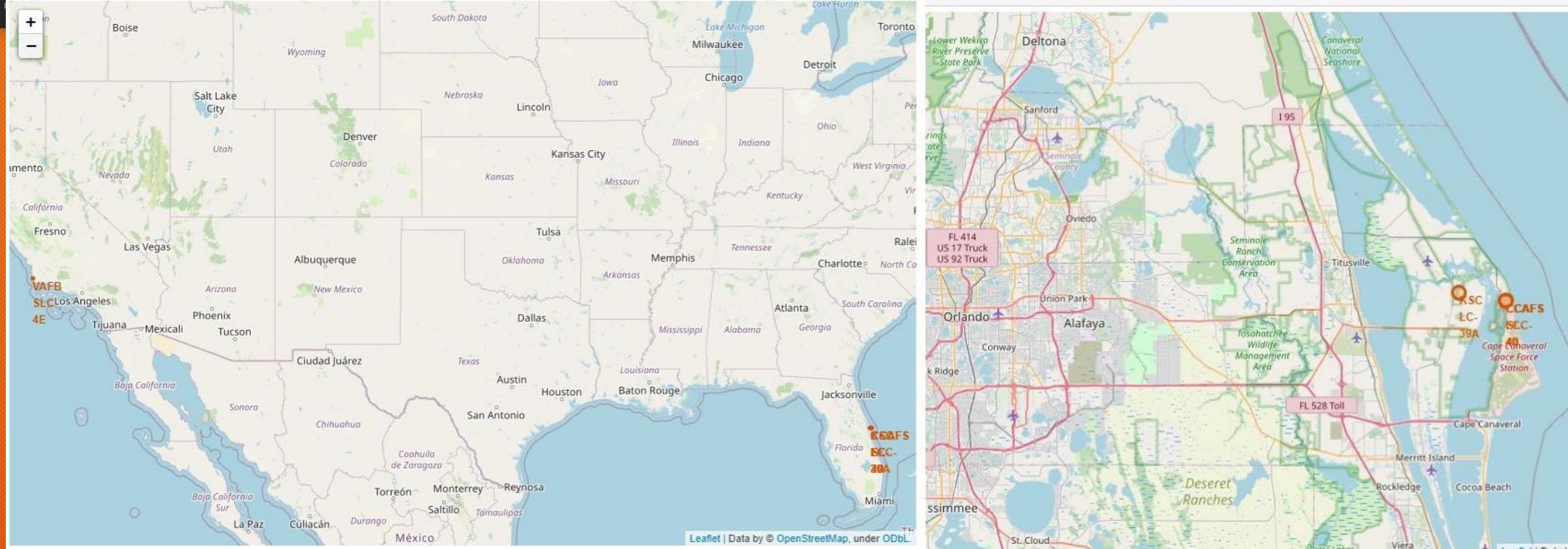
There are two types of successful landing outcomes: drone ship and ground pad landings.

There were 8 successful landings in total during this time period

Interactive Map with Folium

Launch Site Locations

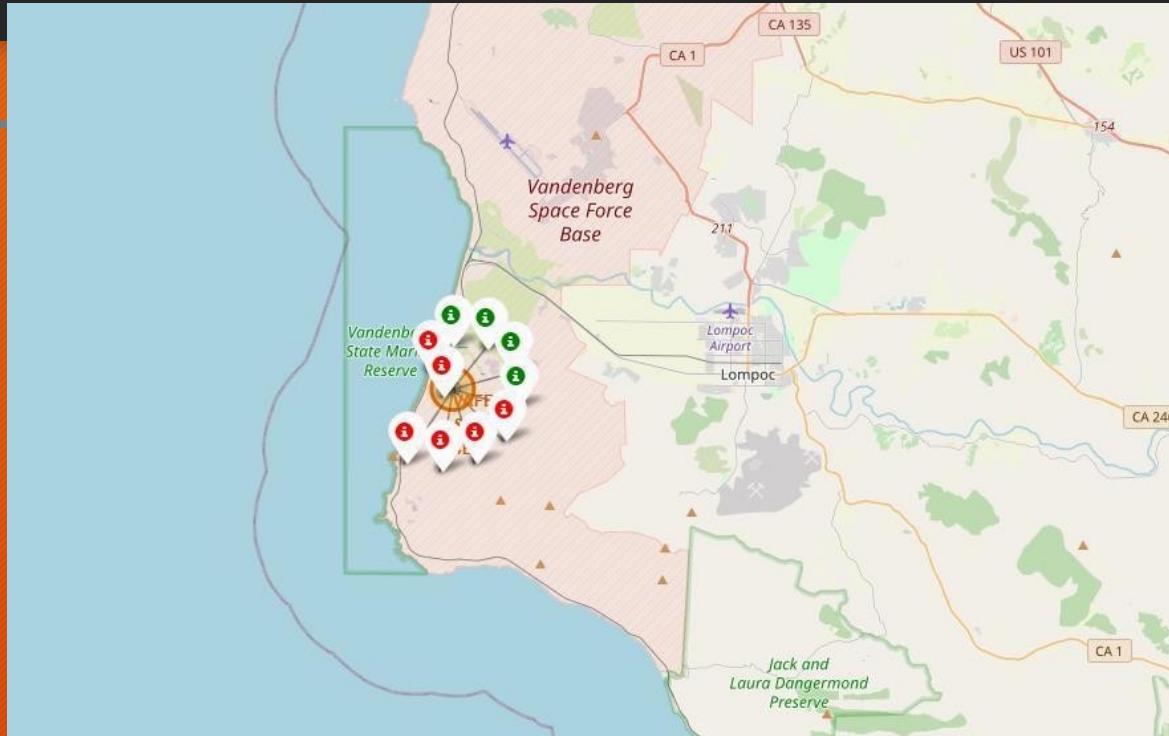
36



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

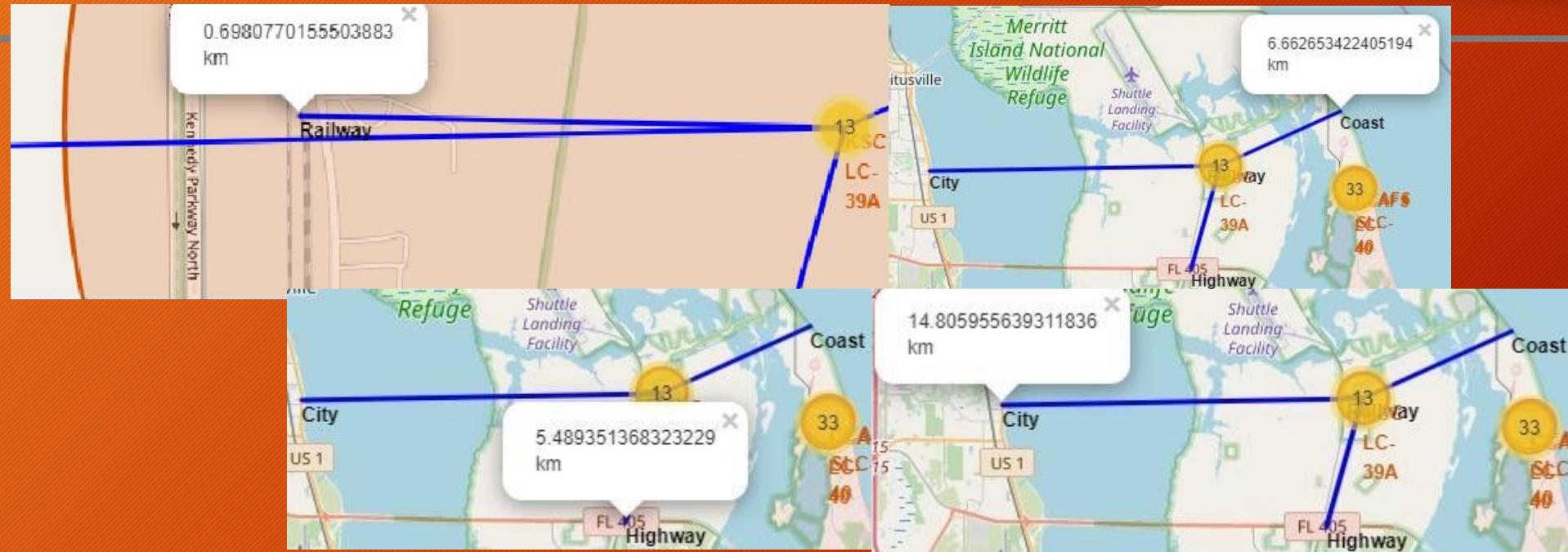
37



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

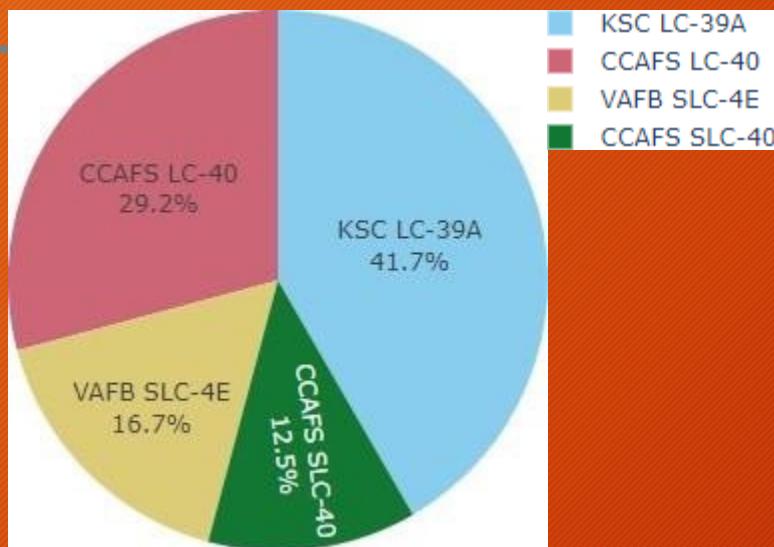
38



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

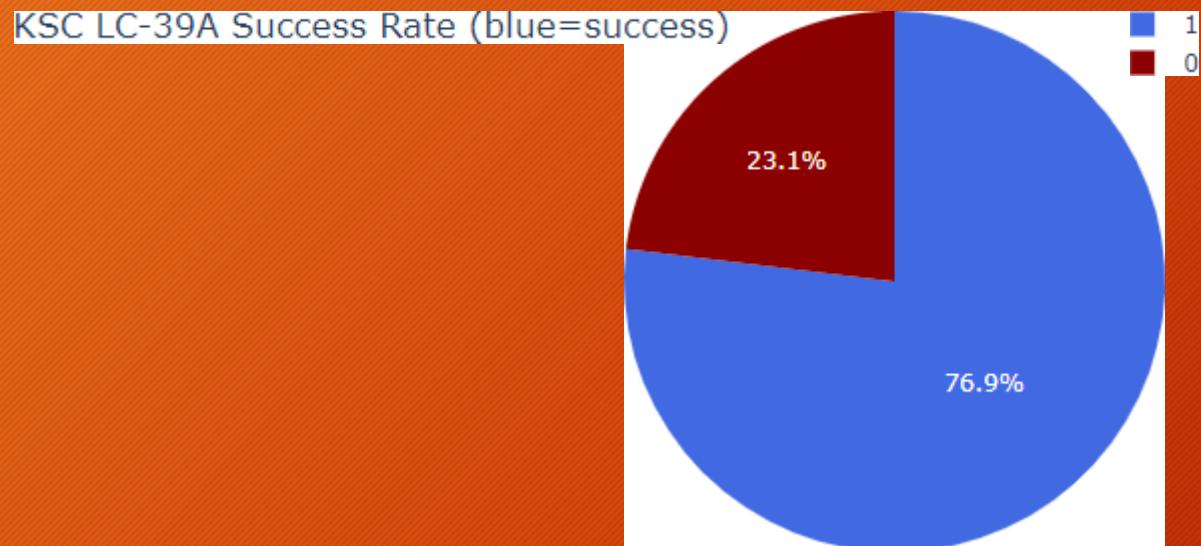
Successful Launches Across Launch Sites



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

41



KSC LC-39A has the highest success rate with 10 successful landings and 3 failed landings.

Payload Mass vs. Success vs. Booster Version Category

42



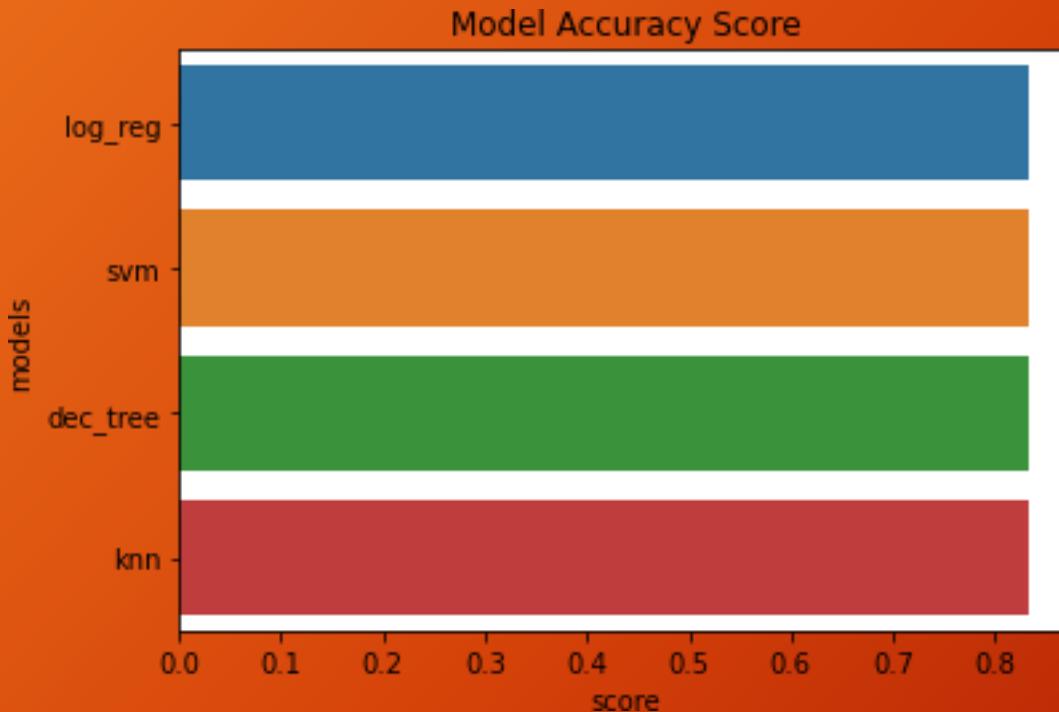
Plotly dashboard has a Payload range selector. However, this is set from 0-10000 instead of the max Payload of 15600. Class indicates 1 for successful landing and 0 for failure. Scatter plot also accounts for booster version category in color and number of launches in point size. In this particular range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

44



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

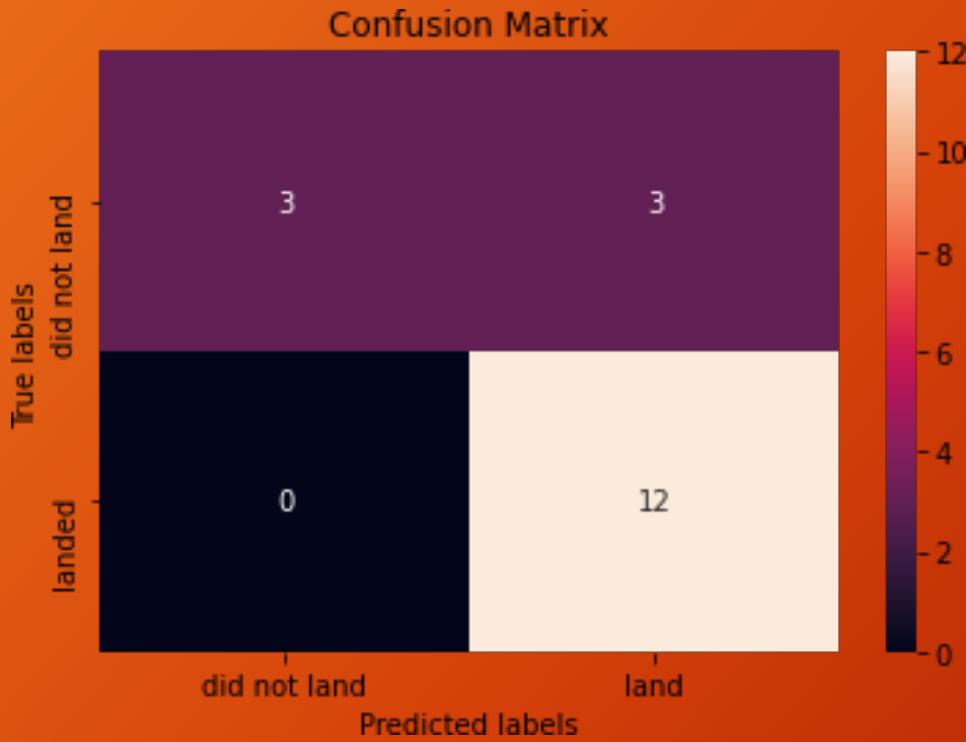
It should be noted that test size is small at only sample size of 18.

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

45



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.

- Our task was to develop a machine learning model for Space Y to help them compete against SpaceX. The model's objective is to predict the successful landing of Stage 1, potentially saving around \$100 million USD. We utilized data from a public SpaceX API and performed web scraping on the SpaceX Wikipedia page. Data labels were created and stored in a DB2 SQL database, and a dashboard was developed for visualization purposes.
- We achieved an accuracy of 83% with our machine learning model. Allon Mask of SpaceY can use this model to predict with reasonable accuracy whether Stage 1 will successfully land before the launch, aiding in the decision-making process regarding the viability of a launch.
- To further improve the model's accuracy, collecting more data would be beneficial. This would help in determining the best machine learning model and enhancing the overall prediction performance.