# Winning Space Race with Data Science

Giru Haran
1 February 2024

# Outline

- Executive Summary

- Introduction

- Methodology

- Results

- Conclusion

- Appendix

# Executive Summary

- Summary of methodologies

    - SpaceX Data Collection using SpaceX API

    - SpaceX Data Collection with Web Scraping

    - SpaceX Data Wrangling

    - SpaceX Exploratory Data Analysis using SQL

    - Space-X EDA DataViz Using Python Pandas and Matplotlib

    - Space-X Launch Sites Analysis with Folium-Interactive Visual Analytics and Ploty Dash

    - SpaceX Machine Learning Landing Prediction

- Summary of all results

    - EDA results

    - Interactive Visual Analytics and Dashboards

    - Predictive Analysis(Classification)

# Introduction



- Project background and context
  SpaceX advertises Falcon 9 rocket launches on its website with a cost of 62 million dollars; other providers cost upward of 165 million dollars each, much of the savings is because SpaceX can reuse the first stage. Therefore, if we can determine if the first stage will land, we can determine the cost of a launch. This information can be used if an alternate company wants to bid against SpaceX for a rocket launch.

- Problems you want to find answers
  In this capstone, we will predict if the Falcon 9 first stage will land successfully using data from Falcon 9 rocket launches advertised on its website.

Section 1

# Methodology

# Methodology

Executive Summary

- Data collection methodology:

  - Describe how data was collected

- Perform data wrangling

  - Describe how data was processed

- Perform exploratory data analysis (EDA) using visualization and SQL

- Perform interactive visual analytics using Folium and Plotly Dash

- Perform predictive analysis using classification models

  - How to build, tune, evaluate classification models

# Data Collection

- Description of how SpaceX Falcon9 data was collected.
    - Data was first collected using SpaceX API (a RESTful API) by making a get request to the SpaceX API. This was done by first defining a series helper functions that would help in the use of the API to extract information using identification numbers in the launch data and then requesting rocket launch data from the SpaceX API URL.

    - Finally, to make the requested JSON results more consistent, the SpaceX launch data was requested and parsed using the GET request and then decoded the response content as a Json result which was then converted into a Pandas data frame.

    - Also performed web scraping to collect Falcon 9 historical launch records from a Wikipedia page titled List of Falcon 9 and Falcon Heavy launches of the launch records are stored in a HTML. Using BeautifulSoup and request Libraries, I extract the Falcon 9 launch HTML table records from the Wikipedia page, Parsed the table and converted it into a Pandas data frame.

# Data Collection – SpaceX API

- Data collected using SpaceX API (a RESTful API) by making a get request to the SpaceX API then requested and parsed the SpaceX launch data using the GET request and decoded the response content as a Json result which was then converted into a Pandas data frame

- Here is the GitHub URL of the completed SpaceX API calls notebook

```python
spacex_url="https://api.spacexdata.com/v4/launches/past"
```

```python
response = requests.get(spacex_url)
```

```python
response_json = response.json()
data = pd.json_normalize(response_json)
```

# Data Collection - Scraping

Performed web scraping to collect Falcon 9 historical launch records from a Wikipedia using BeautifulSoup and request, to extract the Falcon 9 launch records from HTML table of the Wikipedia page, then created a data frame by parsing the launch HTML.

• Here is the GitHub URL of the completed web scraping notebook.

```python
response = requests.get(static_url)
```

```python
# Use BeautifulSoup() to create a BeautifulSoup object
soup = BeautifulSoup(response.content, 'html.parser')
```

```python
html_tables = soup.find_all('table')
```

# Data Wrangling

- After obtaining and creating a Pandas DF from the collected data, data was filtered using the Booster Version column to only keep the Falcon 9 launches, then dealt with the missing data values in the Landing Pad and Payload Mass columns. For the Payload Mass ,missing data values were replaced using mean value of column.

- Also performed some Exploratory Data Analysis (EDA) to find some patterns in the data and determine what would be the label for training supervised models

- Here is the GitHub URL of the completed data wrangling related notebooks.

# EDA with Data Visualization

• Performed data Analysis and Feature Engineering using Pandas and Matplotlib.i.e.

• Exploratory Data Analysis

• Preparing Data Feature Engineering

• Used scatter plots to Visualize the relationship between Flight Number and LaunchSite, Payload and Launch Site, FlightNumber and Orbit type, Payload and Orbittype.

• Used Bar chart to Visualize the relationship between success rate of each orbittype

• Line plot to Visualize the launch success yearly trend.

• Here is the GitHub URL of your completed EDA with data visualization notebook,

# EDA with SQL

- The following SQL queries were performed for EDA

  - Display the names of the unique launch sites in the space mission.

    ```
    %sql SELECT DISTINCT LAUNCH_SITE as "Launch_Sites" FROM SPACEXTBL;
    ```

  - Display 5 records where launch sites begin with the string 'CCA'

    ```
    %sql SELECT * FROM 'SPACEXTBL' WHERE Launch_Site LIKE 'CCA%' LIMIT 5;
    ```

  - Display the total payload mass carried by boosters launched by NASA (CRS)

    ```
    %sql SELECT SUM(PAYLOAD_MASS__KG_) as "Total Payload Mass(Kgs)", Customer FROM 'SPACEXTBL' WHERE Customer = 'NASA (CRS)';
    ```

  - Display average payload mass carried by booster version F9 v1.1

```
%sql SELECT AVG(PAYLOAD_MASS__KG_) as "Payload Mass Kgs", Customer, Booster_Version FROM 'SPACEXTBL' WHERE Booster_Version LIKE 'F9 v1.1%';
```

# EDA with SQL (Cont…)

- List the date when the first successful landing outcome in ground pad was achieved

```
%sql SELECT MIN(DATE) FROM 'SPACEXTBL' WHERE "Landing _Outcome" = "Success (ground pad)";
```

- List the names of the boosters which have success in drone ship and have payload mass greater than 4000 but less than 6000

```
%sql SELECT DISTINCT Booster_Version, Payload FROM SPACEXTBL
    WHERE "Landing_Outcome" = "Success (drone ship)"
        AND PAYLOAD_MASS__KG > 40000 AND AND PAYLOAD_MASS__KG < 60000;
```

- List the total number of successful and failure mission outcomes

```
%sql SELECT "Mission_Outcome", COUNT("Mission_Outcome") as Total FROM SPACEXTBL GROUP BY "Mission_Outcome";
```

- [Here](#) is the GitHub URL of your completed EDA with SQL notebook.

13

# Build an Interactive Map with Folium

- Created folium map to marked all the launch sites, and created map objects such as markers, circles, lines to mark the success or failure of launches for each launch site.

- Created a launch set outcomes (failure=0 or success=1).

- [Here](#) is the GitHub URL of the completed interactive map with Folium map, as an external reference and peer-review purpose.

# Build a Dashboard with Plotly Dash

- Built an interactive dashboard application with Plotly dash by:
  - Adding a Launch Site Drop-down Input Component
  - Adding a callback function to render success-pie-chart based on selected site dropdown
  - Adding a Range Slider to Select Payload
  - Adding a callback function to render the success-payload-scatter-chart scatter plot
- [Here](#) is the GitHub URL of your completed Plotly Dash lab.

# Predictive Analysis (Classification)

Test data using various ML models include SVM, Classification Trees, k nearest neighbors and Logistic Regression;

1. First created an object for each of the algorithms then created a GridSearchCV object and assigned them a set of parameters for each model.

2. GridsearchCV object was created with cv=10, then fit the training data into the GridSearch object for each to Find best Hyperparameter.

3. Output GridSearchCV object for each of the models, then displayed the best parameters using the data attribute best_params_ and the accuracy on the validation data using the data attribute best_score_.

4. Finally using the method score to calculate the accuracy on the test data for each model and plotted a confussion matrix for each using the test and predicted outcomes.

5. The table shows the test data accuracy score for each of the methods comparing them to show which performed best using the test data between SVM, Classification Trees, k nearest neighbors and LogisticRegression;

Here the GitHub URL of the completed predictive analysis lab

# Results

- Exploratory data analysis results

- Interactive analytics demo in screenshots
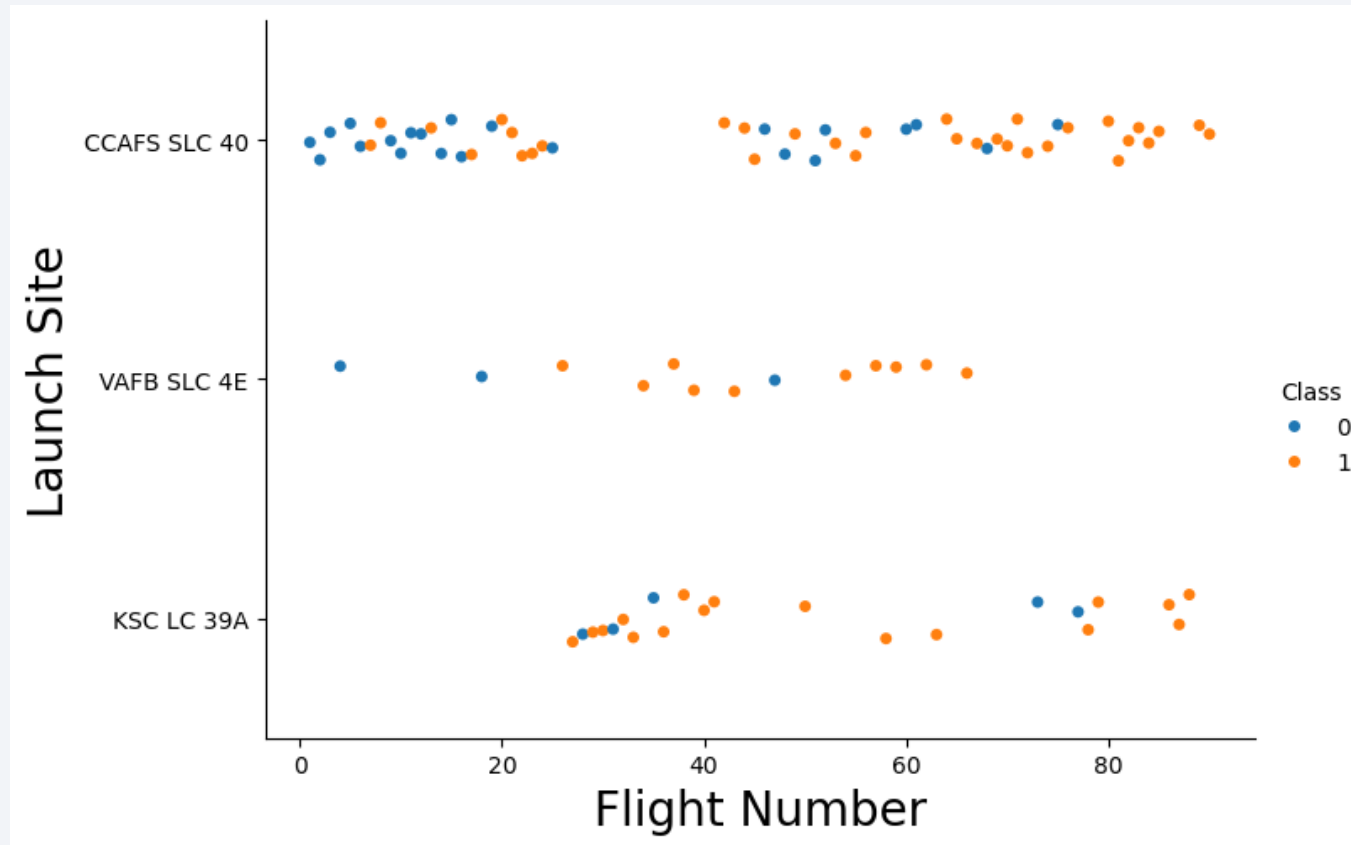
- Predictive analysis results

Refer Github [link](link).

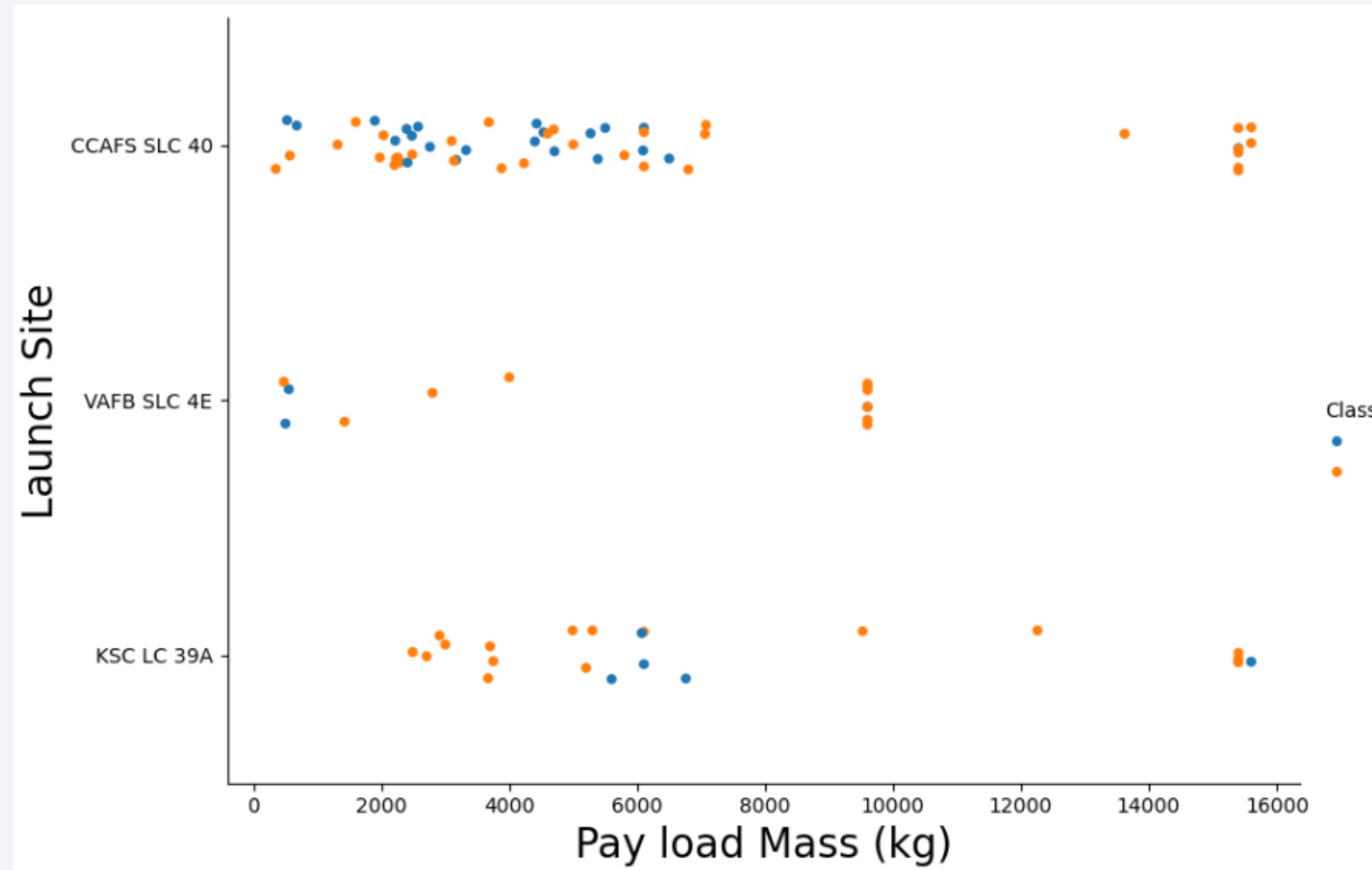Section 2

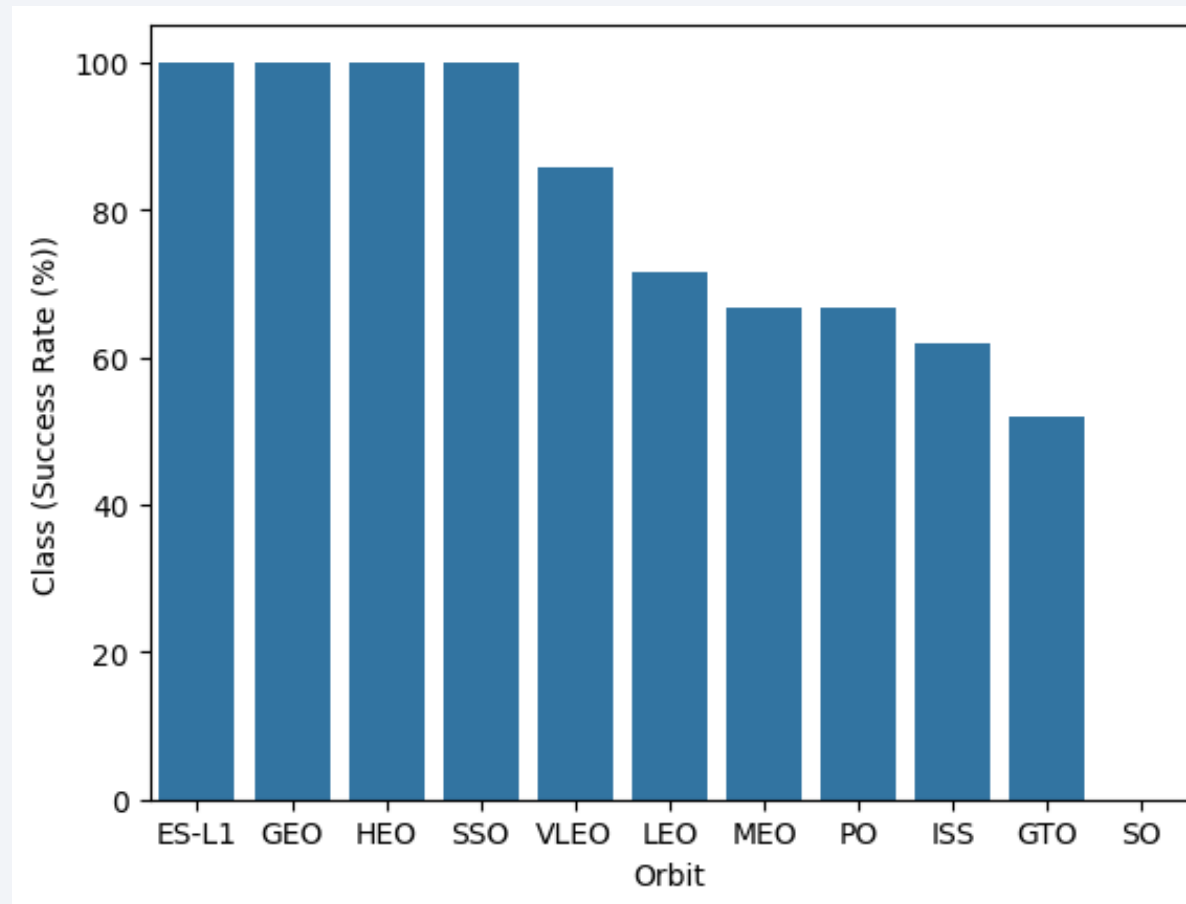# Insights drawn from EDA

# Flight Number vs. Launch Site



It can be inferred that with an increase in the flight number at each of the three launch sites, the success rate also rises. Specifically, for the VAFB SLC 4E launch site, the success rate reaches 100% after the 50th flight. Similarly, both KSC LC 39A and CCAFS SLC 40 achieve a 100% success rate after the 80th flight.
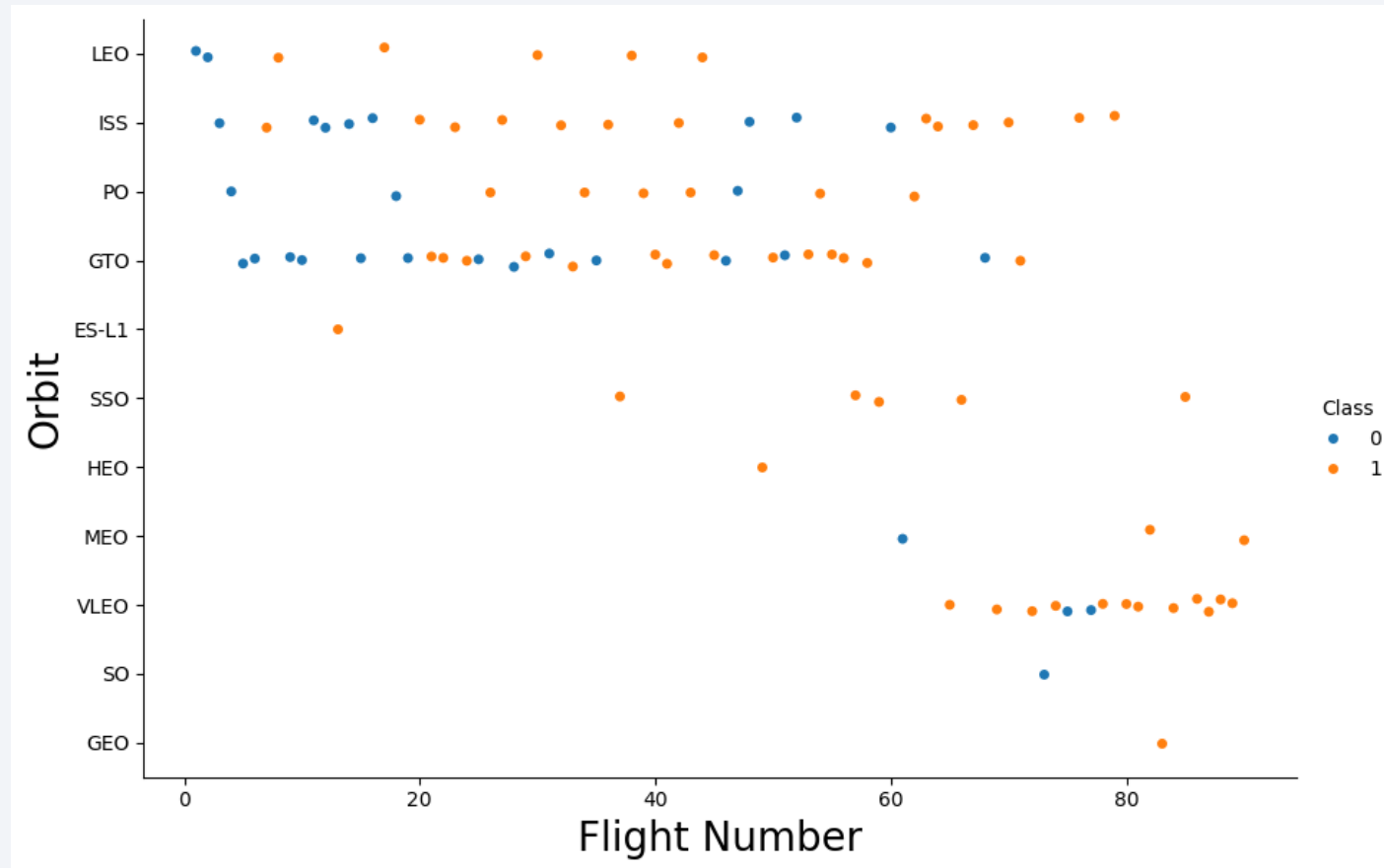
# Payload vs. Launch Site



The VAFB-SLC launch site there are no rockets launched for heavy payload mass(greater than 10000).
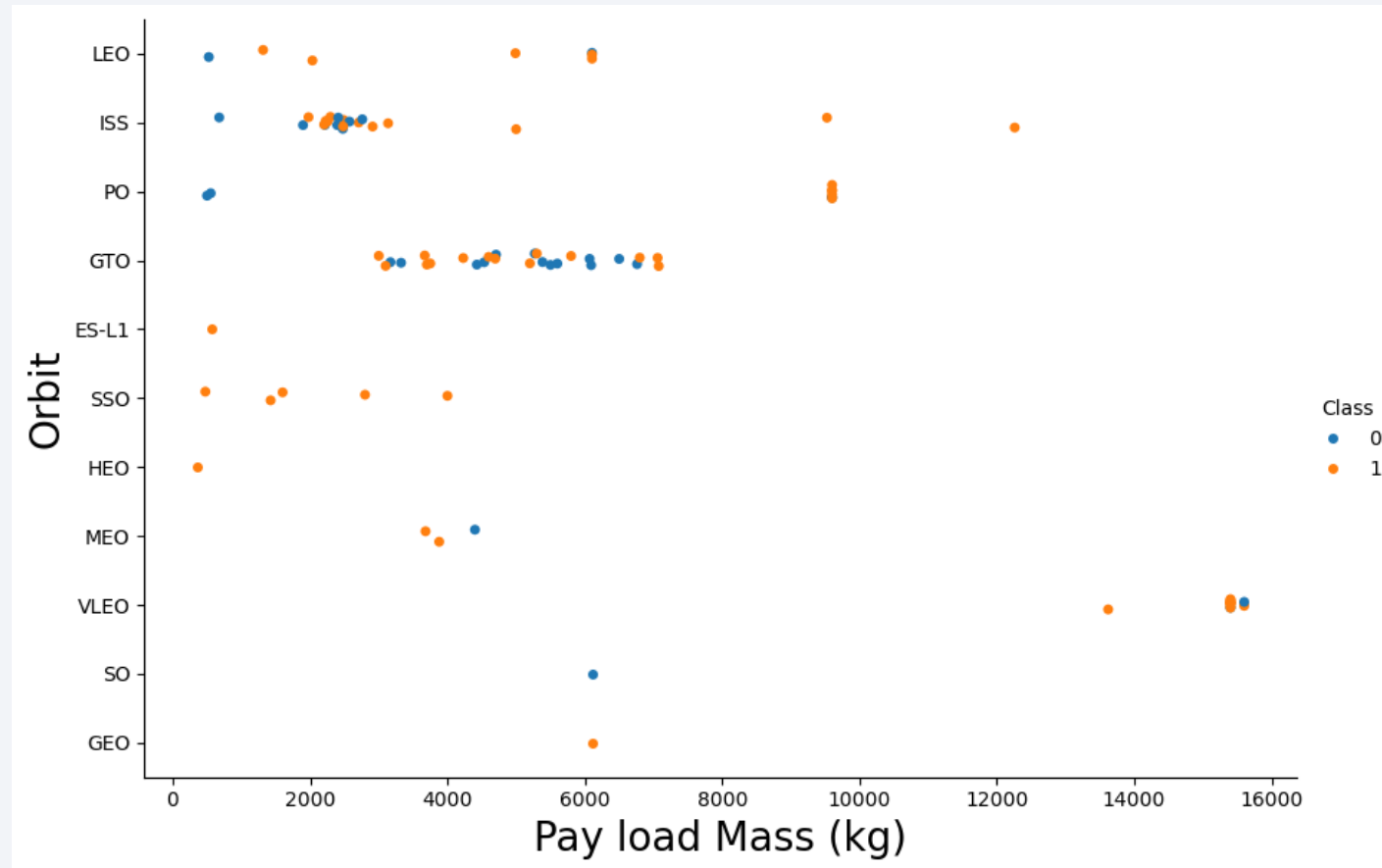
# Success Rate vs. Orbit Type



Orbits ES-L1, GEO, HEO & SSO have the highest success rates at 100%, with SO orbit having the lowest success rate at 0%.

# Flight Number vs. Orbit Type



The LEO orbit the Success appears related to the number of flights; on the other hand, there seems to be no relationship between flight number when in GTO orbit.
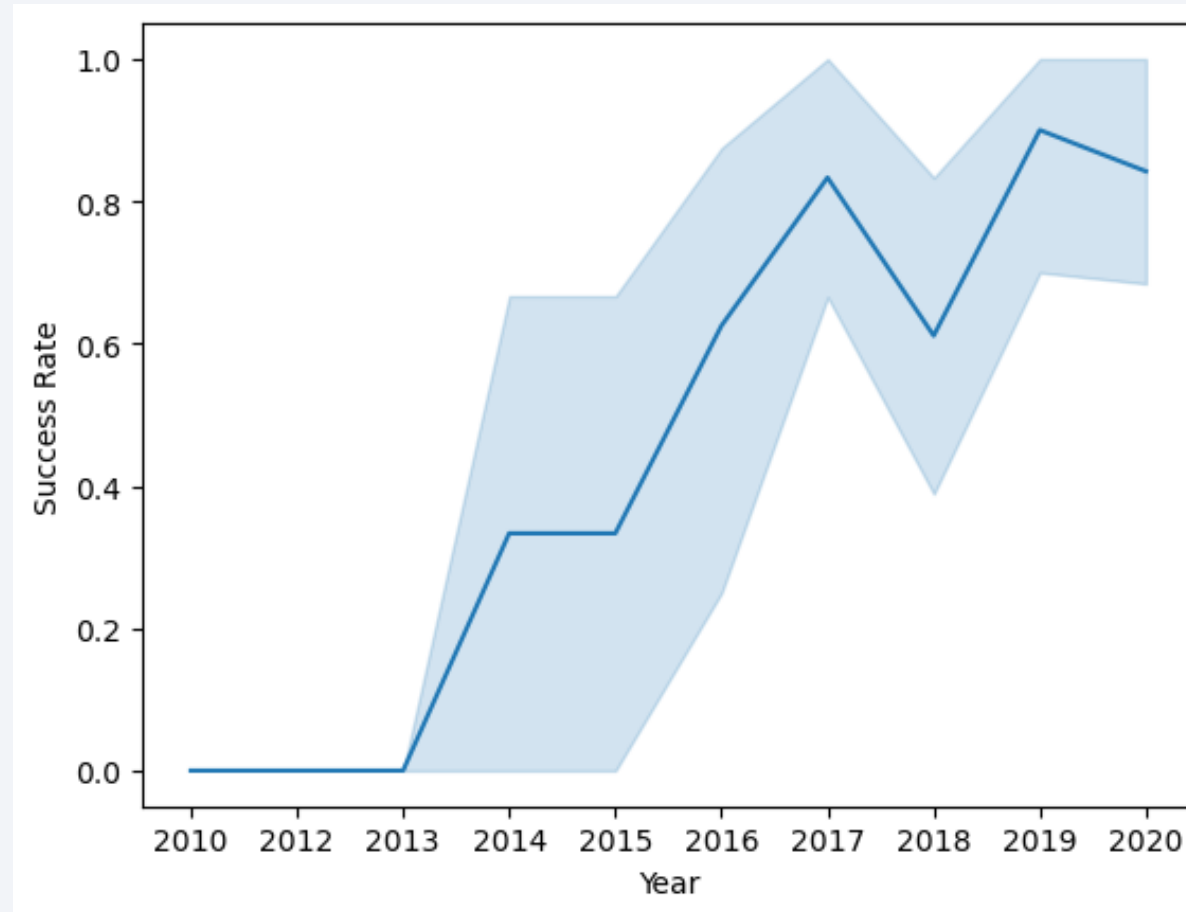
# Payload vs. Orbit Type



With heavy payloads the successful landing or positive landing rate are more for Polar, LEO and ISS. However for GTO we cannot distinguish this well as both positive landing rate and negative landing(unsuccessful mission) are both there here.

# Launch Success Yearly Trend



The success rate since 2013 kept increasing till 2020

# All Launch Site Names

```
%sql SELECT DISTINCT LAUNCH_SITE as "Launch_Sites" FROM SPACEXTBL;
```

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

**Launch_Sites**

CCAFS LC-40

VAFB SLC-4E

KSC LC-39A                    4 Launch sites of East and West Coast Region

CCAFS SLC-40

# Launch Site Names Begin with 'CCA'

```sql
%sql SELECT * FROM 'SPACEXTBL' WHERE Launch_Site LIKE 'CCA%' LIMIT 5;
```

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

| 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 |

# Total Payload Mass

Display the total payload mass carried by boosters launched by NASA (CRS)

```
%sql SELECT SUM(PAYLOAD_MASS__KG_) as "Total Payload Mass(Kgs)", Customer FROM 'SPACEXTBL' WHERE Customer = 'NASA (CRS)';
```

* sqlite:///my_data1.db
Done.

| Total Payload Mass(Kgs) | Customer |
|---|---|
| 45596 | NASA (CRS) |

# Average Payload Mass by F9 v1.1

```
%sql SELECT AVG(PAYLOAD_MASS__KG_) as "Payload Mass Kgs", Customer, Booster_Version FROM 'SPACEXTBL' WHERE Booster_Version
```

* sqlite:///my_data1.db
Done.

| Payload Mass Kgs | Customer | Booster_Version |
|---|---|---|
| 2534.6666666666665 | MDA | F9 v1.1 B1003 |

Payload of 2534 KG is the mean PAYLOAD_MASS__KG_ of all Falcon 9 Missions

# First Successful Ground Landing Date

```
%sql SELECT MIN(Date) FROM 'SPACEXTBL' WHERE "Landing_Outcome" = "Success (ground pad)";
```

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

**MIN(Date)**

2015-12-22

December 2015, Falcon 9 became the first rocket to land propulsively after delivering a payload into orbit. This reusability results in significantly reduced launch costs, as the cost of the first stage constitutes the majority of the cost of a new rocket

# Successful Drone Ship Landing with Payload between 4000 and 6000

```sql
%sql SELECT DISTINCT Booster_Version, Payload FROM SPACEXTBL WHERE "Landing_Outcome" = "Success (drone ship)" AND PAYLOAD_M/
```

* sqlite:///my_data1.db
Done.

| Booster_Version | Payload |
|---|---|
| F9 FT B1022 | JCSAT-14 |
| F9 FT B1026 | JCSAT-16 |
| F9 FT B1021.2 | SES-10 |
| F9 FT B1031.2 | SES-11 / EchoStar 105 |

# Total Number of Successful and Failure Mission Outcomes

```
%sql SELECT "Mission_Outcome", COUNT("Mission_Outcome") as Total FROM SPACEXTBL GROUP BY "Mission_Outcome";
```

* sqlite:///my_data1.db
Done.

| Mission_Outcome | Total |
|---|---|
| Failure (in flight) | 1 |
| Success | 98 |
| Success | 1 |
| Success (payload status unclear) | 1 |

SpaceX Falcon 9 recorded extremely high success rate in main missions about 99% success.

# Boosters Carried Maximum Payload

```
%sql SELECT "Booster_Version",Payload, "PAYLOAD_MASS__KG_" FROM SPACEXTBL WHERE "PAYLOAD_MASS__KG_" = (SELECT MAX("PAYLOAD_I
```

* sqlite:///my_data1.db
Done.

| Booster_Version | Payload | PAYLOAD_MASS__KG_ |
|---|---|---|
| F9 B5 B1048.4 | Starlink 1 v1.0, SpaceX CRS-19 | 15600 |
| F9 B5 B1049.4 | Starlink 2 v1.0, Crew Dragon in-flight abort test | 15600 |
| F9 B5 B1051.3 | Starlink 3 v1.0, Starlink 4 v1.0 | 15600 |
| F9 B5 B1056.4 | Starlink 4 v1.0, SpaceX CRS-20 | 15600 |
| F9 B5 B1048.5 | Starlink 5 v1.0, Starlink 6 v1.0 | 15600 |
| F9 B5 B1051.4 | Starlink 6 v1.0, Crew Dragon Demo-2 | 15600 |
| F9 B5 B1049.5 | Starlink 7 v1.0, Starlink 8 v1.0 | 15600 |
| F9 B5 B1060.2 | Starlink 11 v1.0, Starlink 12 v1.0 | 15600 |
| F9 B5 B1058.3 | Starlink 12 v1.0, Starlink 13 v1.0 | 15600 |
| F9 B5 B1051.6 | Starlink 13 v1.0, Starlink 14 v1.0 | 15600 |
| F9 B5 B1060.3 | Starlink 14 v1.0, GPS III-04 | 15600 |
| F9 B5 B1049.7 | Starlink 15 v1.0, SpaceX CRS-21 | 15600 |

Total 12 version of booster capable of carrying payload up to 15600 KG

# 2015 Launch Records

**Note: SQLLite does not support monthnames. So you need to use substr(Date, 6,2) as month to get the months and substr(Date,0,5)='2015' for year.**

```
%sql SELECT substr(Date, 6, 2) as Month,"Booster_Version", "Launch_Site", Payload, "PAYLOAD_MASS__KG_", "Mission_Outcome",
```

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

| Month | Booster_Version | Launch_Site | Payload | PAYLOAD_MASS__KG_ | Mission_Outcome | Landing_Outcome |
|---|---|---|---|---|---|---|
| 01 | F9 v1.1 B1012 | CCAFS LC-40 | SpaceX CRS-5 | 2395 | Success | Failure (drone ship) |
| 02 | F9 v1.1 B1013 | CCAFS LC-40 | DSCOVR | 570 | Success | Controlled (ocean) |
| 03 | F9 v1.1 B1014 | CCAFS LC-40 | ABS-3A Eutelsat 115 West B | 4159 | Success | No attempt |
| 04 | F9 v1.1 B1015 | CCAFS LC-40 | SpaceX CRS-6 | 1898 | Success | Failure (drone ship) |
| 04 | F9 v1.1 B1016 | CCAFS LC-40 | Turkmen 52 / MonacoSAT | 4707 | Success | No attempt |
| 06 | F9 v1.1 B1018 | CCAFS LC-40 | SpaceX CRS-7 | 1952 | Failure (in flight) | Precluded (drone ship) |
| 12 | F9 FT B1019 | CCAFS LC-40 | OG2 Mission 2 11 Orbcomm-OG2 satellites | 2034 | Success | Success (ground pad) |

# Rank Landing Outcomes Between 2010-06-04 and 2017-03-20

```
%sql SELECT Landing_Outcome, COUNT(*) as OutcomeCount FROM SPACEXTBL WHERE Date BETWEEN '2010-06-04' AND '2017-03-20' GROUP
```

* sqlite:///my_data1.db
Done.

| Landing_Outcome | OutcomeCount |
| --- | --- |
| No attempt | 10 |
| Success (drone ship) | 5 |
| Failure (drone ship) | 5 |
| Success (ground pad) | 3 |
| Controlled (ocean) | 3 |
| Uncontrolled (ocean) | 2 |
| Failure (parachute) | 2 |
| Precluded (drone ship) | 1 |

Cumulatively 10 success Landing outcome can be observed including drone ship and ground pad.

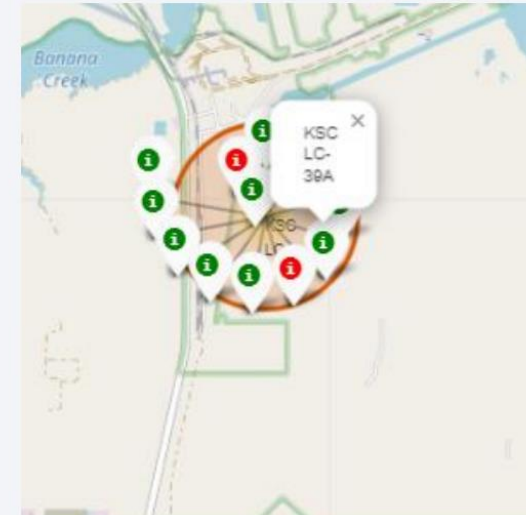Section 3

# Launch Sites
# Proximities Analysis

# Markers of all launch sites on global map



Main Launch site region include Eastern Coast (Florida) and West Coast (California)

# Launch outcomes for each site on the map With Color Markers







Eastern coast (Florida) Launch site KSC LC-39A has relatively high success rates compared to CCAFS SLC-40 & CCAFS LC-40
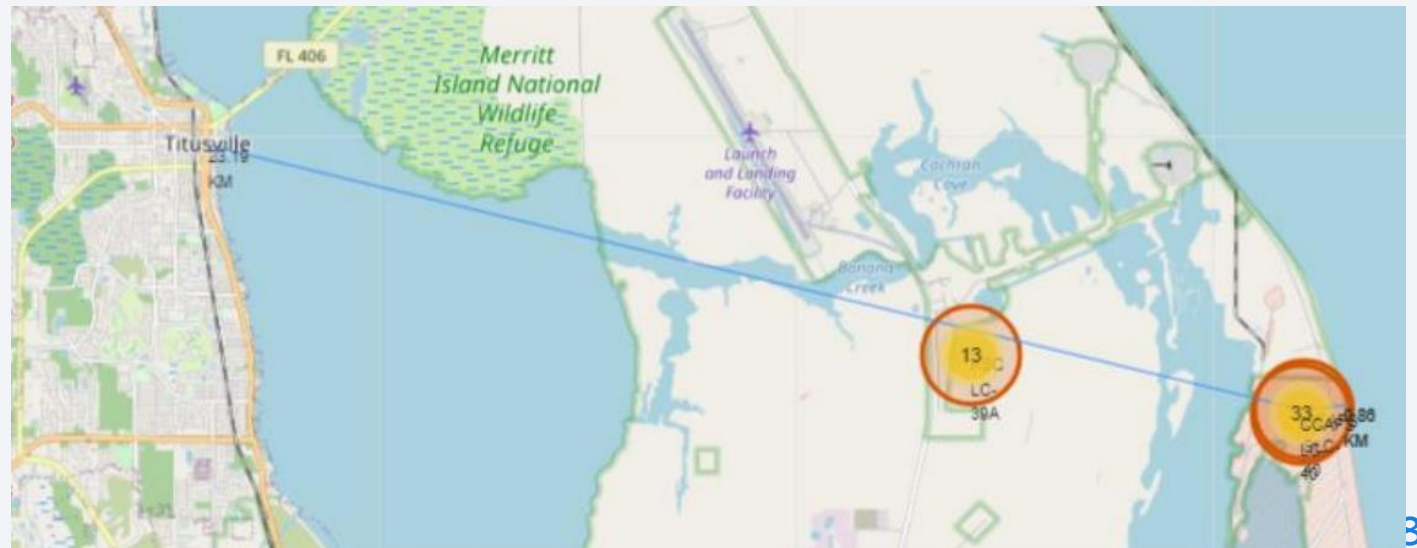


In the West Coast (California) Launch site VAFB SLC-4E has relatively lower success rates 4/10 compared to KSC LC-39A launch site in the Eastern Coast of Florida.

# Distances between a launch site to its proximities



Launch site CCAFS SLC-40 proximity to coastline is 0.86km

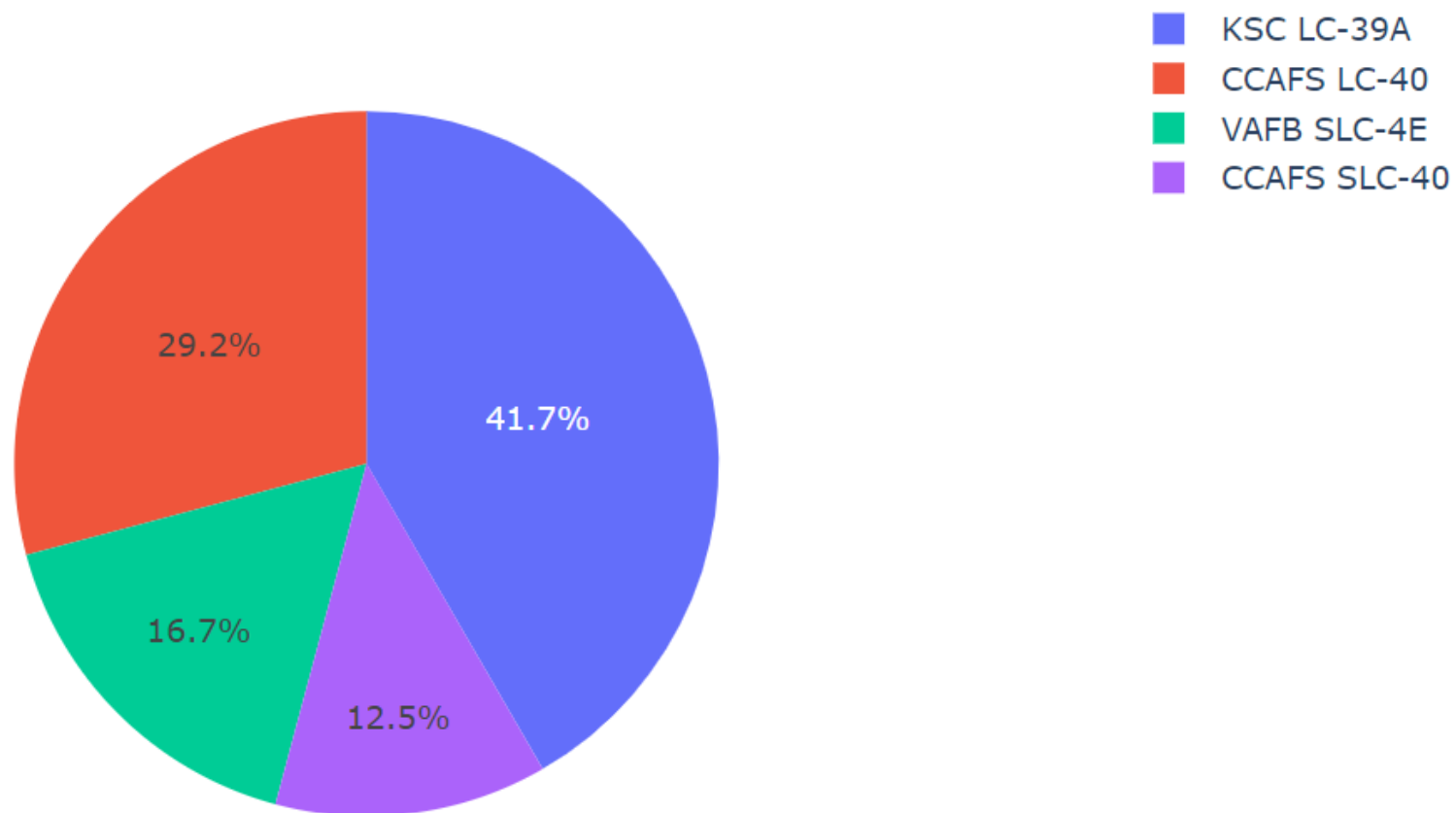Launch site CCAFS SLC-40 closest to highway (Washington Avenue) is 23.19km

# Build a Dashboard with Plotly Dash

# SpaceX Launch Records Dashboard

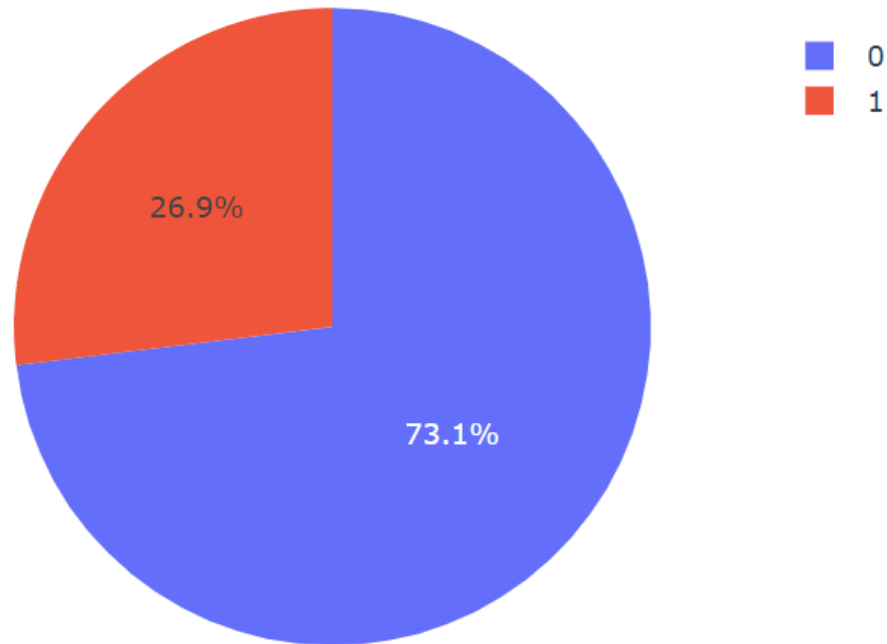## Success Count for all launch sites

Launch site KSC LC-39A has the
highest launch success rate at 42%
followed by CCAFS LC-40 at 29%,
VAFB SLC-4E at 17% and lastly launch
site CCAFS SLC-40 with a success rate
of 13%

**Legend:**
- KSC LC-39A
- CCAFS LC-40
- VAFB SLC-4E
- CCAFS SLC-40

- 41.7% — KSC LC-39A
- 29.2% — CCAFS LC-40
- 16.7% — VAFB SLC-4E
- 12.5% — CCAFS SLC-40
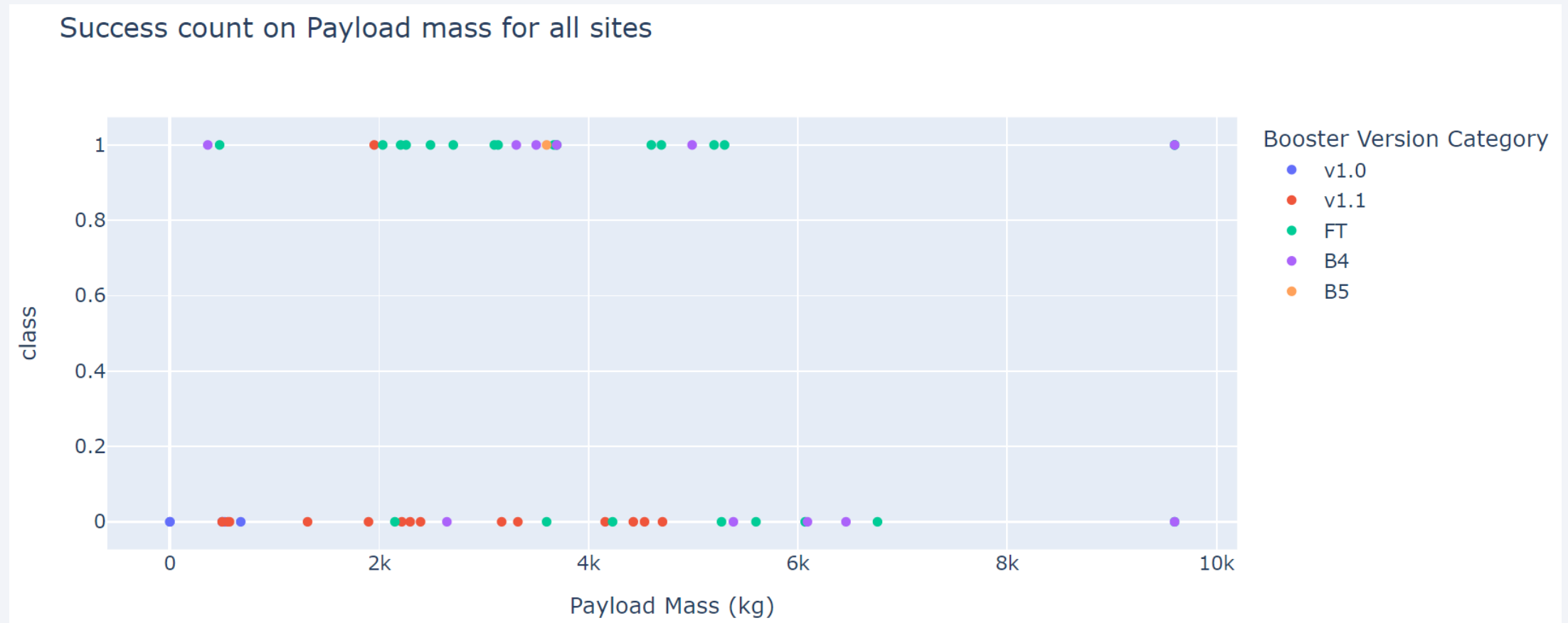
# 2ⁿᵈ Highest Launch Site Record

CCAFS LC-40

Total Success Launches for site CCAFS LC-40



Launch site CCAFS LC-40 had the 2nd highest success ratio of 73% success against 27% failed launches

# Payload vs. Launch Outcome scatter plot for all sites



For Launch site CCAFS LC-40 the booster version FT has the largest success rate from a payload mass of >2000kg

Section 5

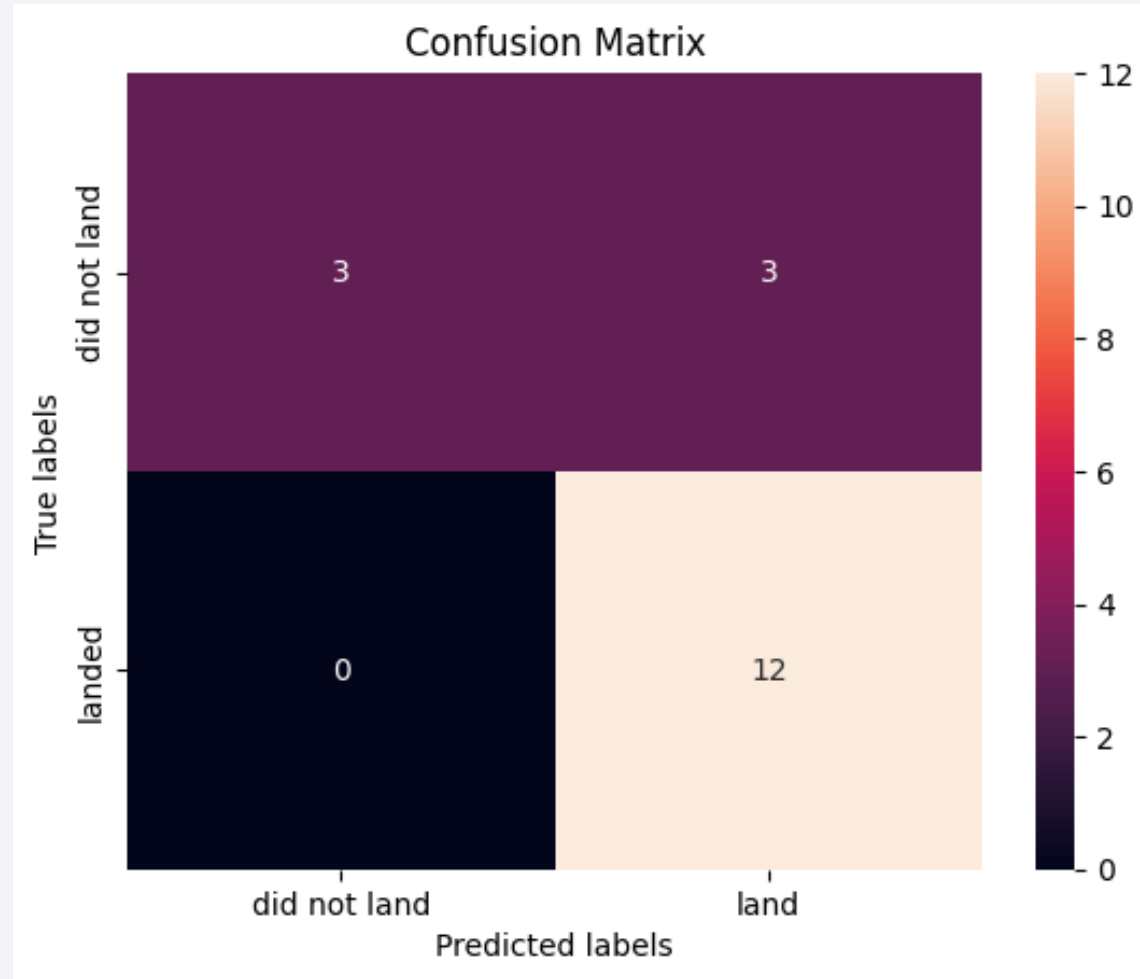# Predictive Analysis (Classification)

# Classification Accuracy

- All the methods **perform equally** on the test data: i.e. They all have the same accuracy of **0.833333** on the test Data.

| Method | Test Data Accuracy |
|---|---|
| Logistic_Reg | 0.833333 |
| SVM | 0.833333 |
| Decision Tree | 0.833333 |
| KNN | 0.833333 |

# Confusion Matrix



All the 4 classification model had the same confusion matrixes and were able equally distinguish between the different classes. The major problem is false positives for all the models.

# Conclusions

• Different launch sites have different success rates. CCAFS LC-40, has a success rate of 60 %,while KSC LC-39A and VAFB SLC 4E has a success rate of 77%.

• We can deduce that, as the flight number increases in each of the 3 launching sites, so does the success rate. For instance, the success rate for the VAFB SLC 4E launch site is 100% after the Flight number 50. Both KSC LC 39A and CCAFS SLC 40 have a 100% success rates after 80thflight

• If you observe Payload Vs. Launch Site scatter point chart you will find for the VAFB-SLC launch site there are no rockets launched for heavy payload mass(greater than 10000).

• Orbits ES-L1, GEO, HEO & SSO have the highest success rates at 100%, with SO orbit having the lowest success rate at ~50%. Orbit SO has 0% success rate.

• LEO orbit the Success appears related to the number of flights; on the other hand, there seems to be no relationship between flight number when in GTO orbitConclusions57

• With heavy payloads the successful landing or positive landing rate are more for Polar, LEO and ISS. However, for GTO we cannot distinguish this well as both positive landing rate and negative landing(unsuccessful mission) are both there here

• The success rate since 2013 kept increasing till 2020.

Thank you!