

The background is a dark blue space-themed illustration. In the center is a teal-colored Earth showing the continents of Europe, Africa, and Asia. Surrounding the Earth are three concentric circular orbits. Several teal-colored circles of varying sizes are placed on these orbits, representing satellites or celestial bodies. The background is also filled with numerous white stars of different sizes and shapes, some appearing as simple dots and others as multi-pointed stars.

# Data science for space industry

Aditya Naik

4<sup>th</sup> September 2023

# Presentation outline

- Executive summary
- Introduction
- Methodology
- Results
- Conclusion
- Appendix



# Executive summary

- Business Understanding: a binary Classification challenge- predicting whether Falcon9 rockets can land successfully or not
- Data Understanding: revealed need for multiple data sources, therefore compiled a dataset from,
  - SpaceX Rest API
  - Webscraping the Falcon rocket Wikipedia page
- Data Preparation: used several Python data manipulation and visualisation libraries, along with SQL to manipulate the data and ensure a clean meaningful dataset stored in a database
- Exploratory Data Analysis: carried out using Pandas and Matplotlib, and Seaborn to identify feature variables
- Key outputs:
  - Developed an interactive leaflet map with Folium, and an interactive dashboard with Plotly Dash
  - Predictive model for Falcon 9 Landings using K-Nearest neighbour
- Outstanding:
  - Evaluation with stakeholders
  - Model deployment
  - Further refinement

# Introduction

- Business Understanding:
- Space travel is an expensive exercise, with each launch costing 62m USD by SpaceX, and almost three times as much from other vendors.
- Predicting whether a SpaceX Falcon9 rocket will land is crucial to deciding a vendor for a rocket launch
- A successful prediction would result in savings of up to tens of millions of dollars if the decision is correct
- Question: Can publicly available data of historic SpaceX launches be used to develop a Machine Learning model to predict whether Falcon9 launches land successfully or not?

# Methodology



## Executive Summary

- Data collection methodology:
  - Compiling a dataset from,
    - SpaceX Rest API
    - Complemented with data from webscraping SpaceX Wikipedia page
      - Using BeautifulSoup method
- Perform data wrangling
  - Payload mass had five missing values: replaced with the mean
- 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



- SpaceX REST API used to gather some of the data
  - Gathered data includes Booster version, longitude and latitudes of launch sites, payload, cores, flight number, and date
  - Filtered out other Falcon versions to leave only Falcon9
- Webscraping of the Falcon9 Wikipedia page used to extract critical tables for the final dataframe
  - Data points under this included: Flight number, Date, Time, Booster Version, Launch site, Payload, Payload mass, Orbit, Customer, Launch Outcome, and Booster Landing
  - Methods used: Get requests, BeautifulSoup, and html.parser

# Data collection – SpaceX API

- After steps in flowchart on right, then filtered out other Falcon versions to leave only Falcon9

1st API call to rockets: For Booster Version names

2nd API Call to launchpads: For launch sites along with their longitudes & latitudes

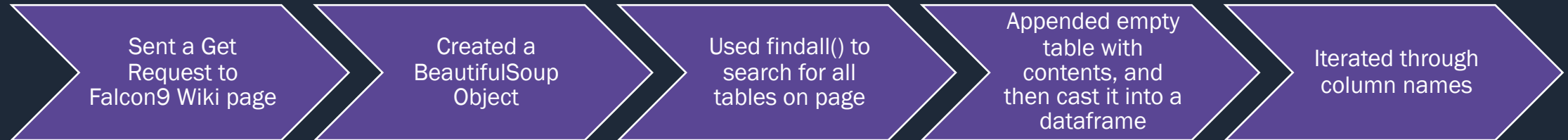
3rd API Call to Payloads: for payload masses and orbits

4th API Call to Cores: for landing outcome, info on cores, gridfins, legs, and landing pad

Data parsed in a Dictionary with column names as keys, and data as values



# Data collection - Scrapping





# Data wrangling

Checked for missing data under each variable and expressed as a %

Checked data types of each variable using method dtypes to prevent issues later in analysis and modeling

Conducted value\_counts() on key variables

Used “Outcome” variable to add new column “Class” depending on whether successful landing or not



# EDA with data visualization



- Used Pandas, Matplotlib and Seaborn to plot
  - Scatterplot of Flight number, Payload mass, and Class
    - To assess whether progressive flights succeed or not, and effect of payload mass on safe return
  - Scatterplot of Flight number, Launch site, and Class
    - To check for striking data points on launch site, and flight number and their effect on class
  - Bar chart showing success rate of each type of orbit
    - To see if any orbits have outstanding outcomes compared to others
  - Scatterplot of Orbit, Flight number, and Class
    - To see if any orbits have striking relationships with flight number and class
  - Scatterplot of Orbit, Payload mass, and Class
    - Assess relationship between Orbit, payload mass and success rate
  - Line plot of mean annual success rate
    - To show whether the program has been improving with time

# EDA with SQL



- A query to identify unique site used for launch by SpaceX
- Showed five records from launch sites with names containing string “CCA”
- Displayed total payload mass carried by boosters launched by NASA (CRS)
- Showed the mean payload mass carried by booster version F9 v1.1
- Listed the date when the first successful landing outcome in ground pad was achieved.
- List the names of the boosters which have success in drone ship and have payload mass greater than 4000 but less than 6000
- List the total number of successful and failure mission outcomes
- Used a subquery to list the names of the Booster versions which have carried the maximum payload mass
- Used Substr() to List the records which will display the month names, failure landing outcomes in drone ship ,booster versions, launch site for the months in year 2015.
- Rank the count of landing outcomes (such as Failure (drone ship) or Success (ground pad)) between the date 2010-06-04 and 2017-03-20, in descending order.

# Build an interactive map with Folium



- Created circles and markers for each of the four launch sites
- Added additional markers to each launch site in which
  - Green markers for successful launches
  - Red markers for unsuccessful launches
- Put a marker on nearest coastline, and calculated distance to a launch site using the geometry of the earth (assuming the earth's radius is 6373km)
- Drew a Polyline between the coastline above, and the launch site
- Put markers on, and then calculated distances to the following objects,
  - Highway (very close, for logistical ease)
  - Railroad (very close for logistical ease)
  - City (further away, for public safety)

# Build a dashboard with plotly Dash



- Added,
  - a “Launch Site” Drop-down Input Component
  - a callback function to render ”Success Pie-chart” based on selected site dropdown
  - a Range Slider to select “Payload”
  - a callback function to render the ”Success Payload Scatter-chart” scatter plot
- These will assist the business colleagues to make informed decisions using an interactive platform to gain insights

# Predictive analysis (Classification)



Created a Numpy array for the target Y from "Class" variable

Used StandardScaler to standardize the data of the Features (X)

Did Train-Test split using 80:20

Identified best-performing approach

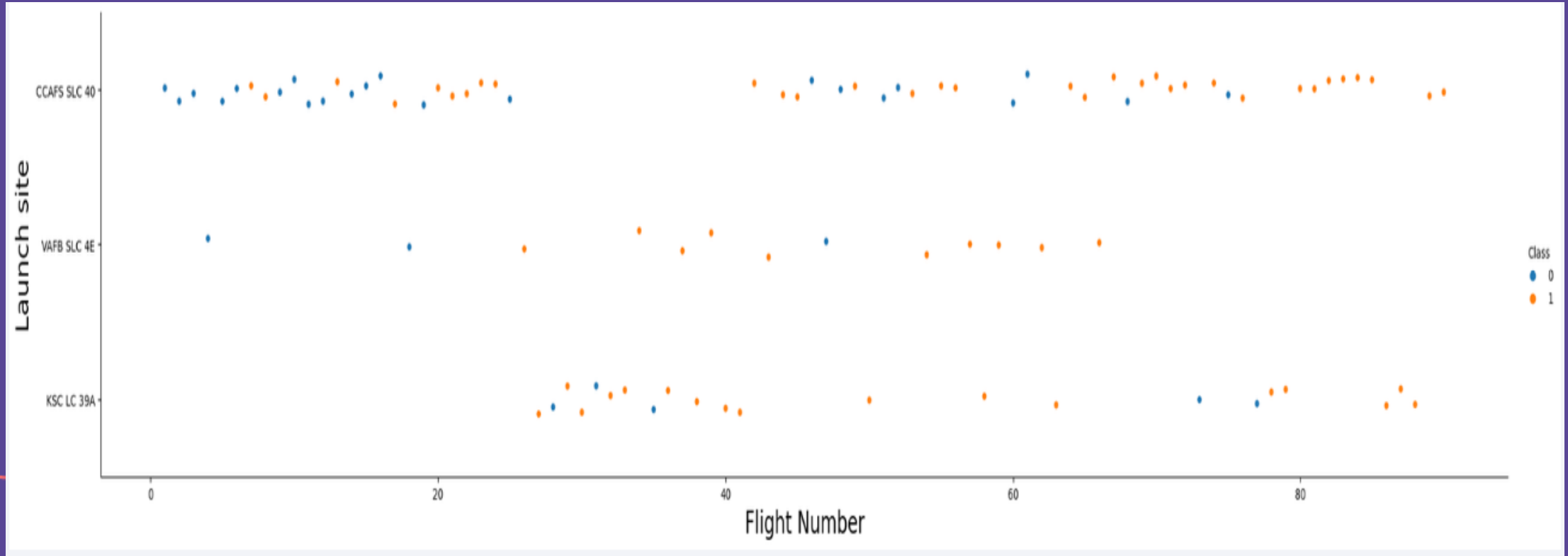
For each approach,  
1. Identified best parameters  
2. Assessed their accuracy using test data  
3. Plotted Confusion Matrix

# Results

- Exploratory data analysis results
- Interactive analytics demo in screenshots
- Predictive analysis results

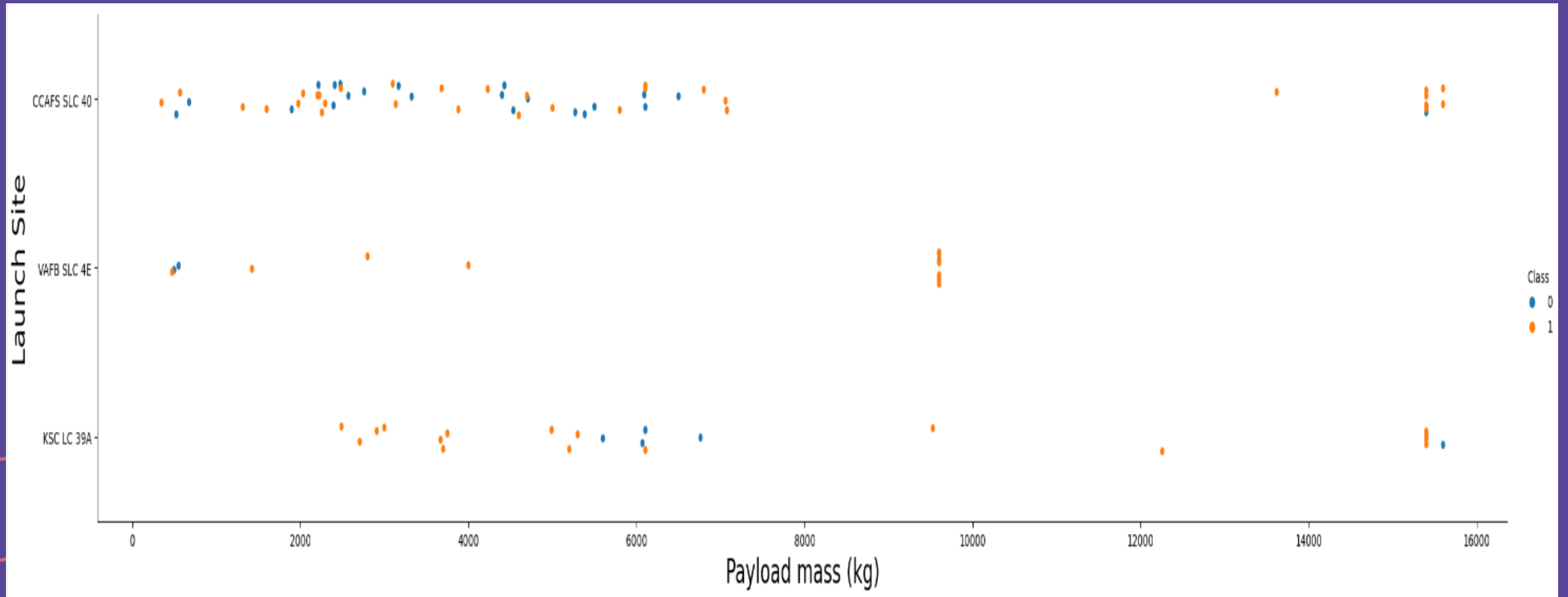


# EDA: Flight number vs Launch Site

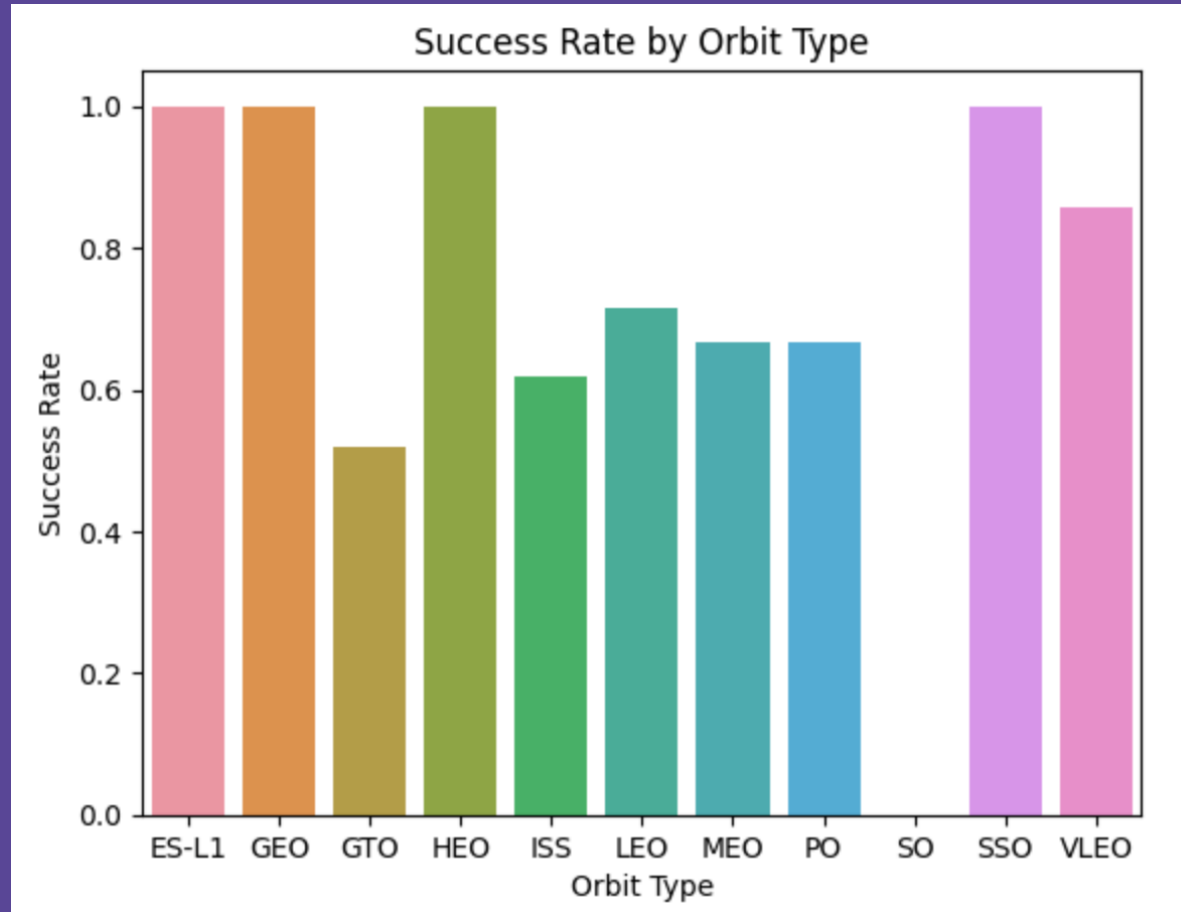




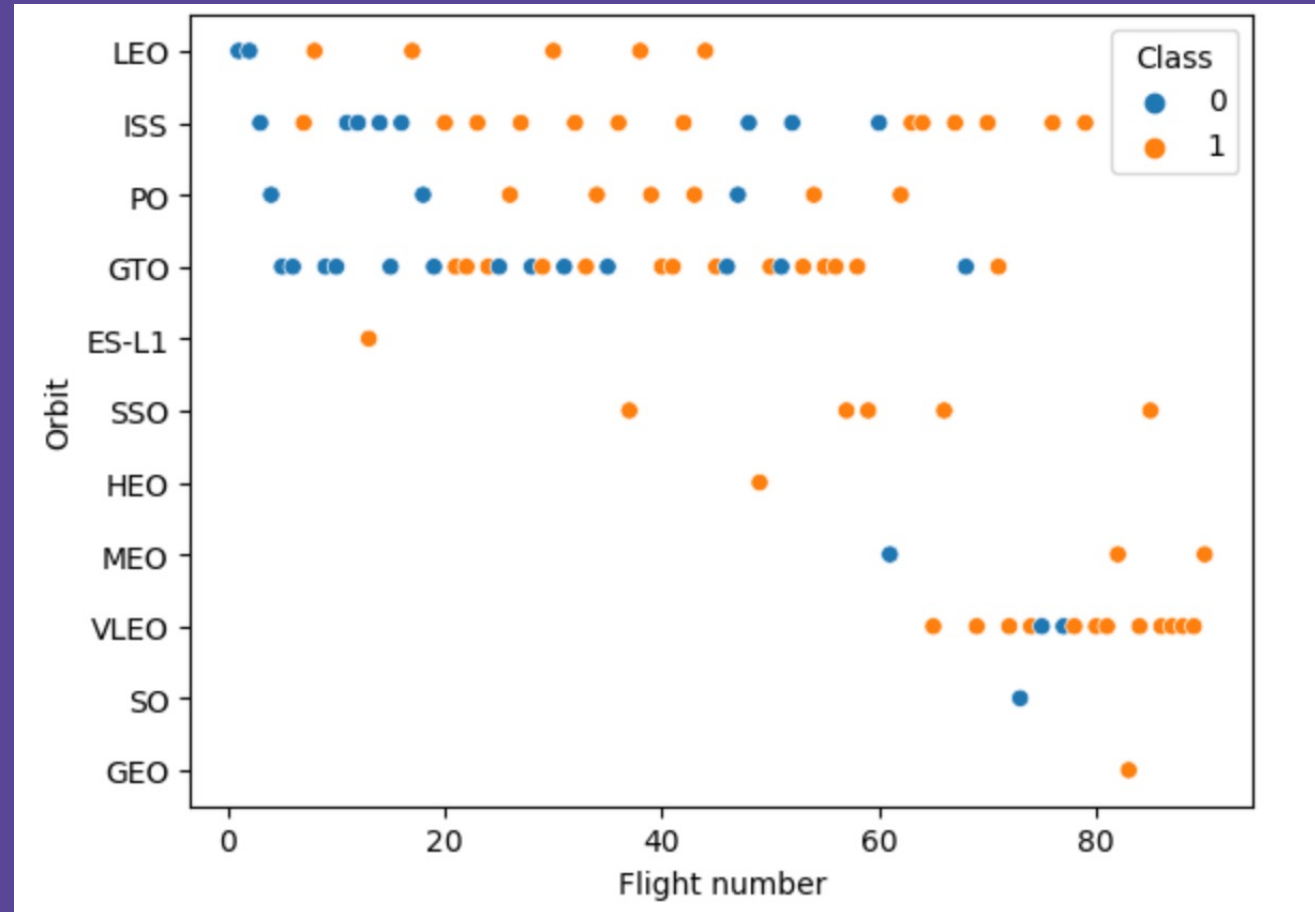
# EDA: Payload vs launch site



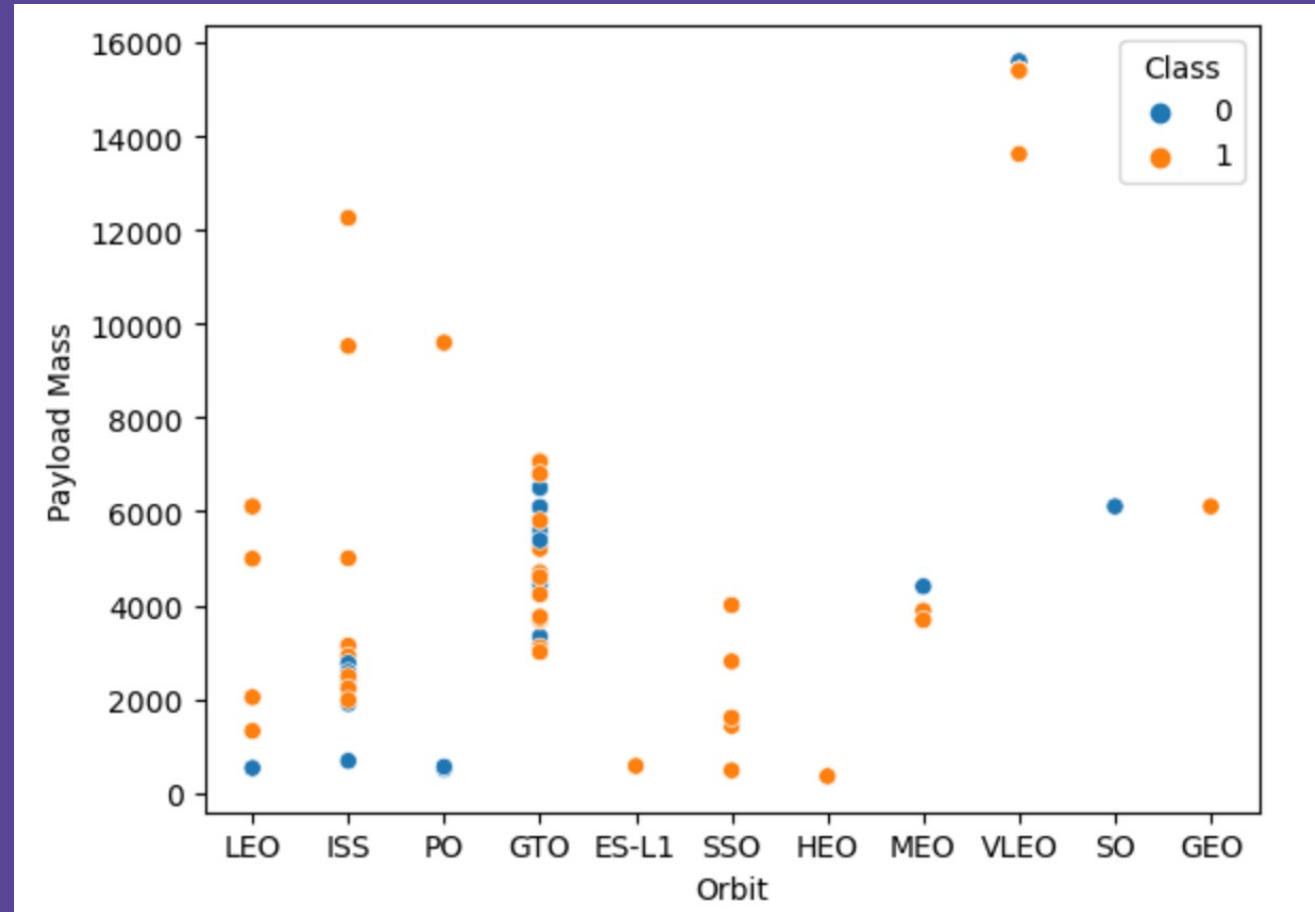
# EDA: Success rate vs orbit type



# EDA: Flight number vs orbit type



# EDA: Payload vs orbit type



# EDA: Launch success yearly trend



# EDA: All launch site names

```
In [16]: 1 %sql select DISTINCT ("Launch_Site") from SPACEXTABLE
```

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

Out[16]:

Launch_Site
-------------

CCAFS LC-40
-------------

VAFB SLC-4E
-------------

KSC LC-39A
------------

CCAFS SLC-40
--------------

# EDA: Launch site names beginning with CCA

```
In [24]: 1 %sql Select * from SPACEXTABLE where "Launch_Site" like "%CCA%" limit 5;
```

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

Out[24]:

Date	Time (UTC)	Booster_Version	Launch_Site	Payload	PAYLOAD_MASS_KG_	Orbit	Customer	Mission_Outcome	Landing_Outcome
2010-04-06	18:45:00	F9 v1.0 B0003	CCAFS LC-40	Dragon Spacecraft Qualification Unit	0	LEO	SpaceX	Success	Failure (parachute)
2010-08-12	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
2012-08-10	00:35:00	F9 v1.0 B0006	CCAFS LC-40	SpaceX CRS-1	500	LEO (ISS)	NASA (CRS)	Success	No attempt
2013-01-03	15:10:00	F9 v1.0 B0007	CCAFS LC-40	SpaceX CRS-2	677	LEO (ISS)	NASA (CRS)	Success	No attempt

# EDA: Total payload mass

```
In [38]: 1 %sql Select SUM("PAYLOAD_MASS_KG_") from SPACEXTABLE where "Customer"="NASA (CRS)";
```

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

```
Out[38]: SUM("PAYLOAD_MASS_KG_")
```

```
45596
```



# EDA: Average payload mass for Falcon 9 v1.1

```
In [39]: 1 %sql select AVG ("PAYLOAD_MASS_KG_") from SPACEXTABLE where "Booster_Version"="F9 v1.1"

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

Out[39]:

AVG ("PAYLOAD_MASS_KG_")
2928.4

# EDA: First successful ground landing date

```
In [50]: 1 %sql Select * from SPACEXTABLE where "Landing_Outcome" like "%success%" order by Date asc limit 1
```

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

```
Done.
```

```
Out[50]:
```

Date	Time (UTC)	Booster_Version	Launch_Site	Payload	PAYLOAD_MASS_KG_	Orbit	Customer	Mission_Outcome	Landing_Outcome
2015-12-22	01:29:00	F9 FT B1019	CCAFS LC-40	OG2 Mission 2 11 Orbcomm-OG2 satellites	2034	LEO	Orbcomm	Success	Success (ground pad)

# EDA: Successful Drone Ship Landing with Payload between 4000 and 6000

```
In [18]: 1 ion") from SPACEXTABLE where "Landing_Outcome"="Success (drone ship)" AND "PAYLOAD_MASS_KG_" between 4000 and 6000
```

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

Done.

```
Out[18]: Booster_Version
```

F9 FT B1022

F9 FT B1026

F9 FT B1021.2

F9 FT B1031.2

# EDA: Total number of Successful and failure missions

```
In [52]: 1 %sql Select "Mission_Outcome", Count("Mission_Outcome") from SPACEXTABLE group by "Mission_Outcome"
```

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

Out[52]:

Mission_Outcome	Count("Mission_Outcome")
Failure (in flight)	1
Success	98
Success	1
Success (payload status unclear)	1

# EDA: Maximum payload for each booster

```
In [63]: 1 %sql Select DISTINCT "Booster_Version", "PAYLOAD_MASS_KG_" from SPACEXTABLE where "PAYLOAD_MASS_KG_" = (Select MAX
          * sqlite:///my_data1.db
          Done.
```

```
Out[63]:
```

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

# EDA: 2015 launch record

```
In [72]: 1 %sql Select "Booster_Version", "Launch_Site" FROM SPACEXTABLE where "Landing_Outcome" = "Failure (drone ship)" AND  
* sqlite:///my_data1.db  
Done.
```

```
Out[72]:
```

Booster_Version	Launch_Site
F9 v1.1 B1012	CCAFS LC-40
F9 v1.1 B1015	CCAFS LC-40

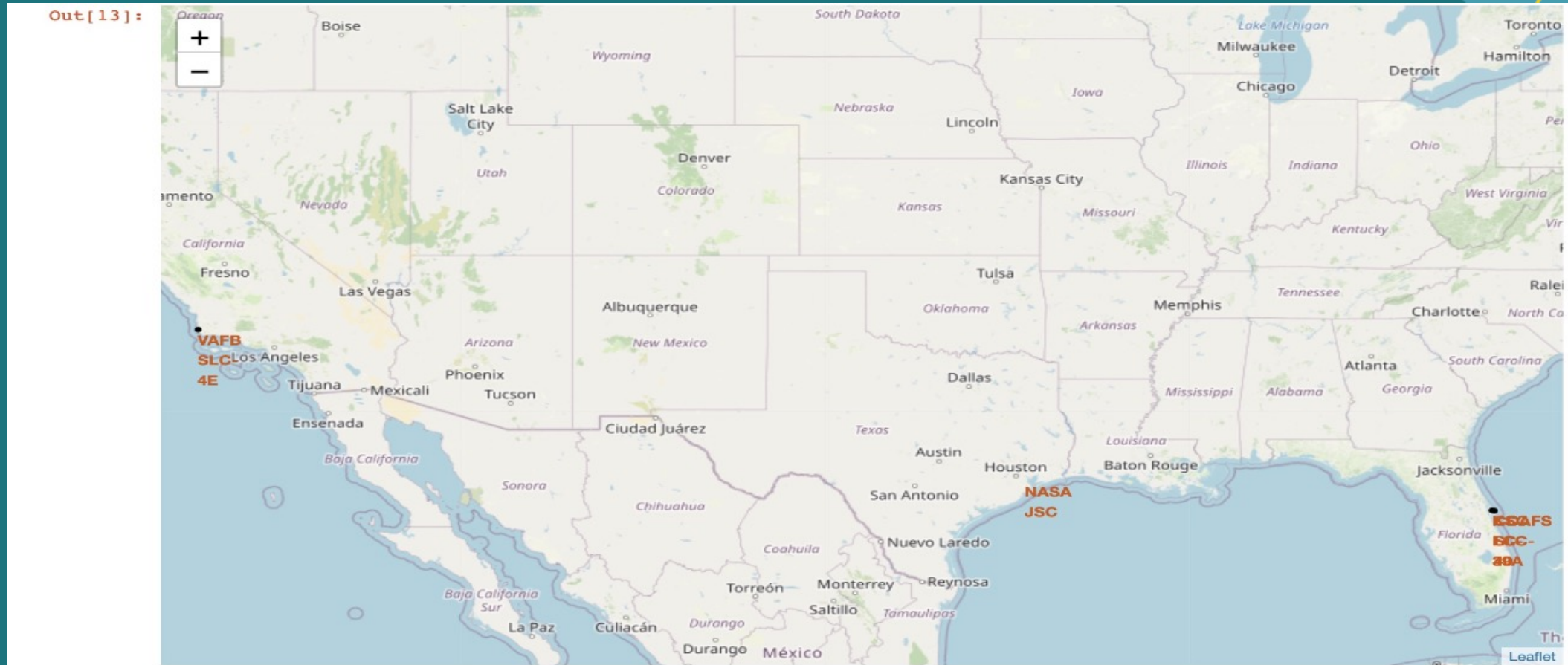
# EDA: Ranking landing outcomes

```
In [65]: 1 %sql Select "Landing_Outcome", count(*) from SPACEXTABLE where Date between '2011-06-04' and '2017-03-20' group by  
* sqlite:///my_data1.db  
Done.
```

```
Out[65]:
```

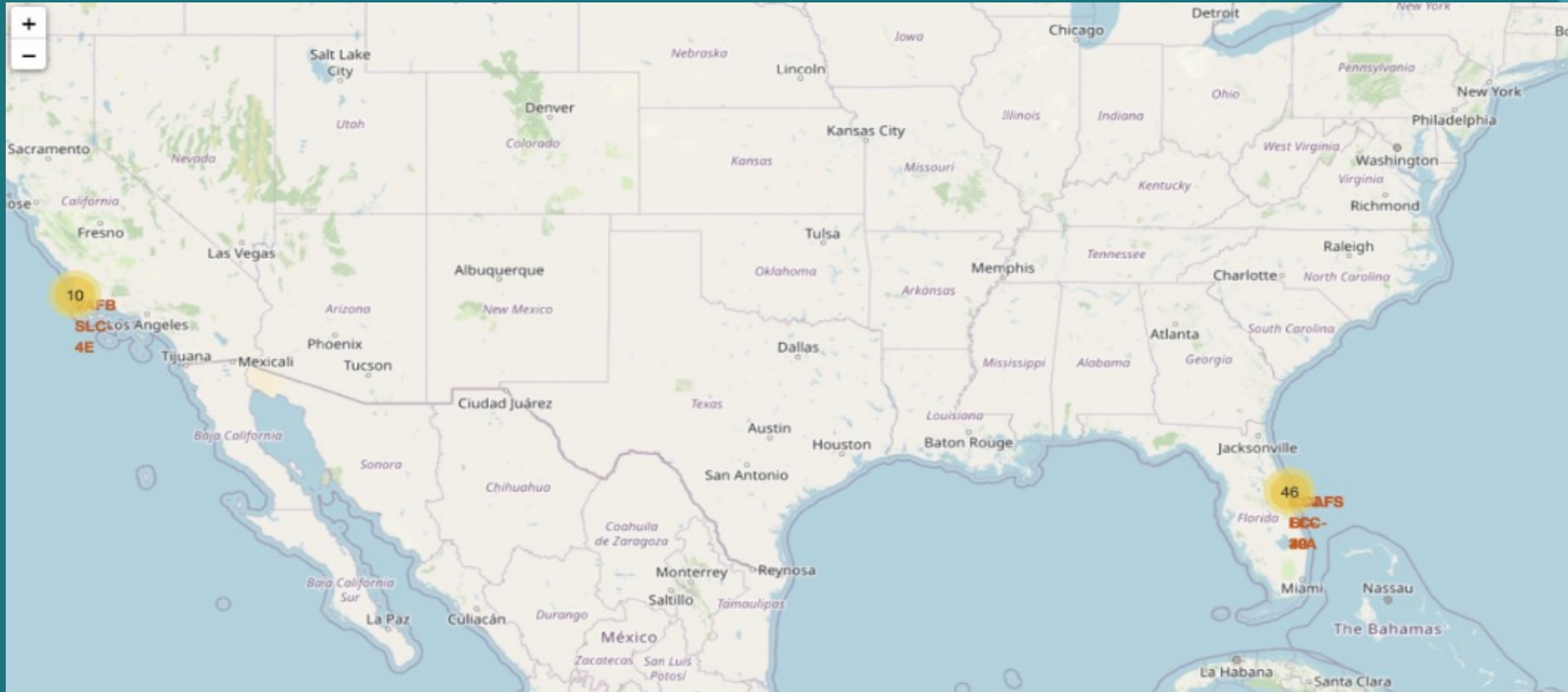
Landing_Outcome	count(*)
No attempt	10
Success (ground pad)	5
Success (drone ship)	5
Failure (drone ship)	5
Controlled (ocean)	3
Uncontrolled (ocean)	2
Precluded (drone ship)	1

# Interactive analysis: Folium leaflet of SpaceX launch sites

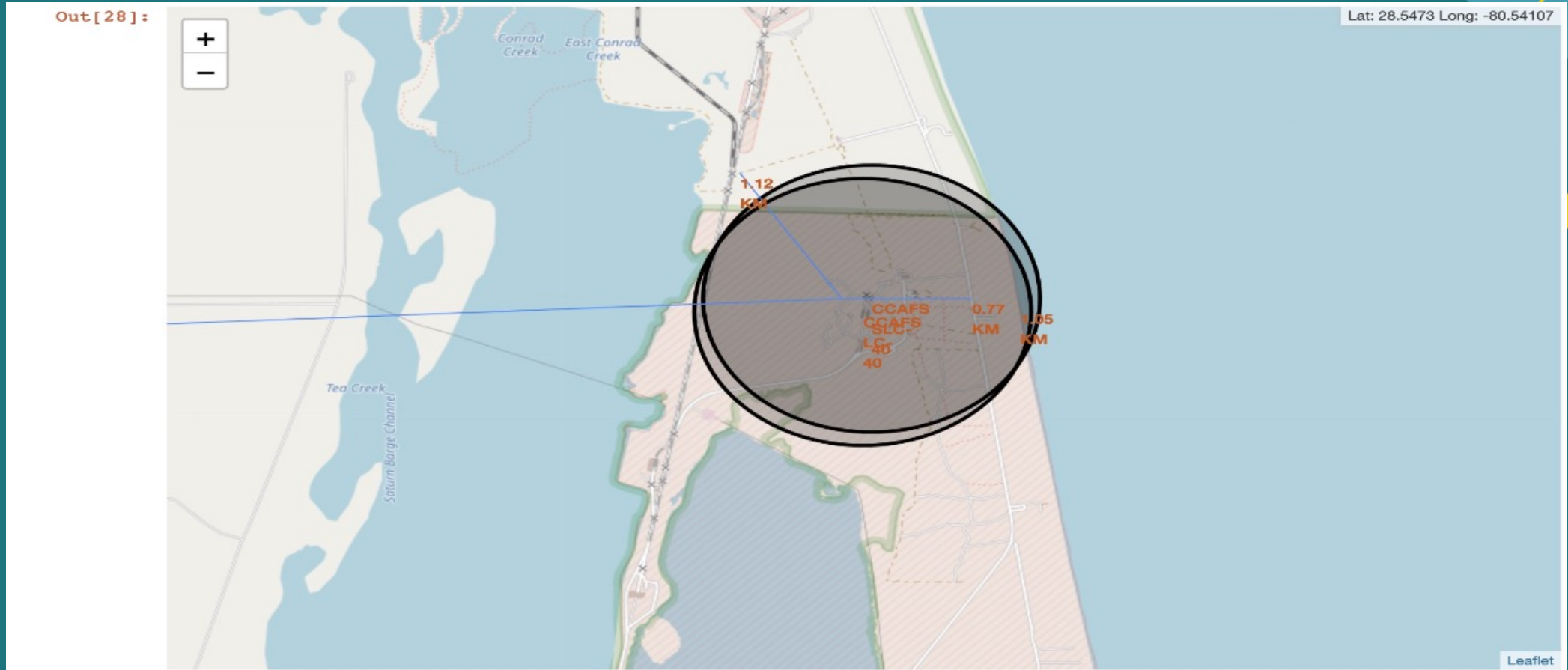




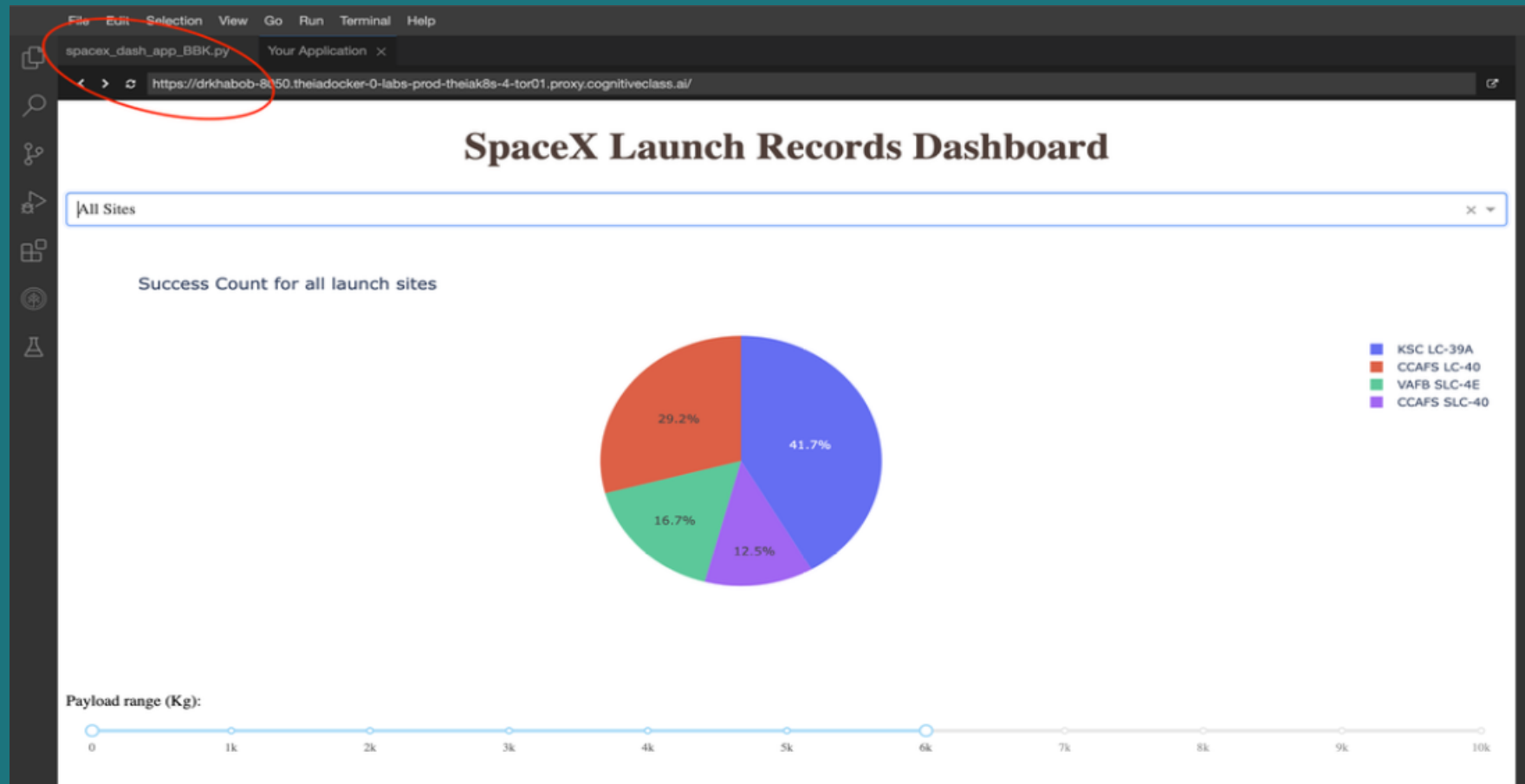
# Interactive analysis: success and failure of SpaceX launches



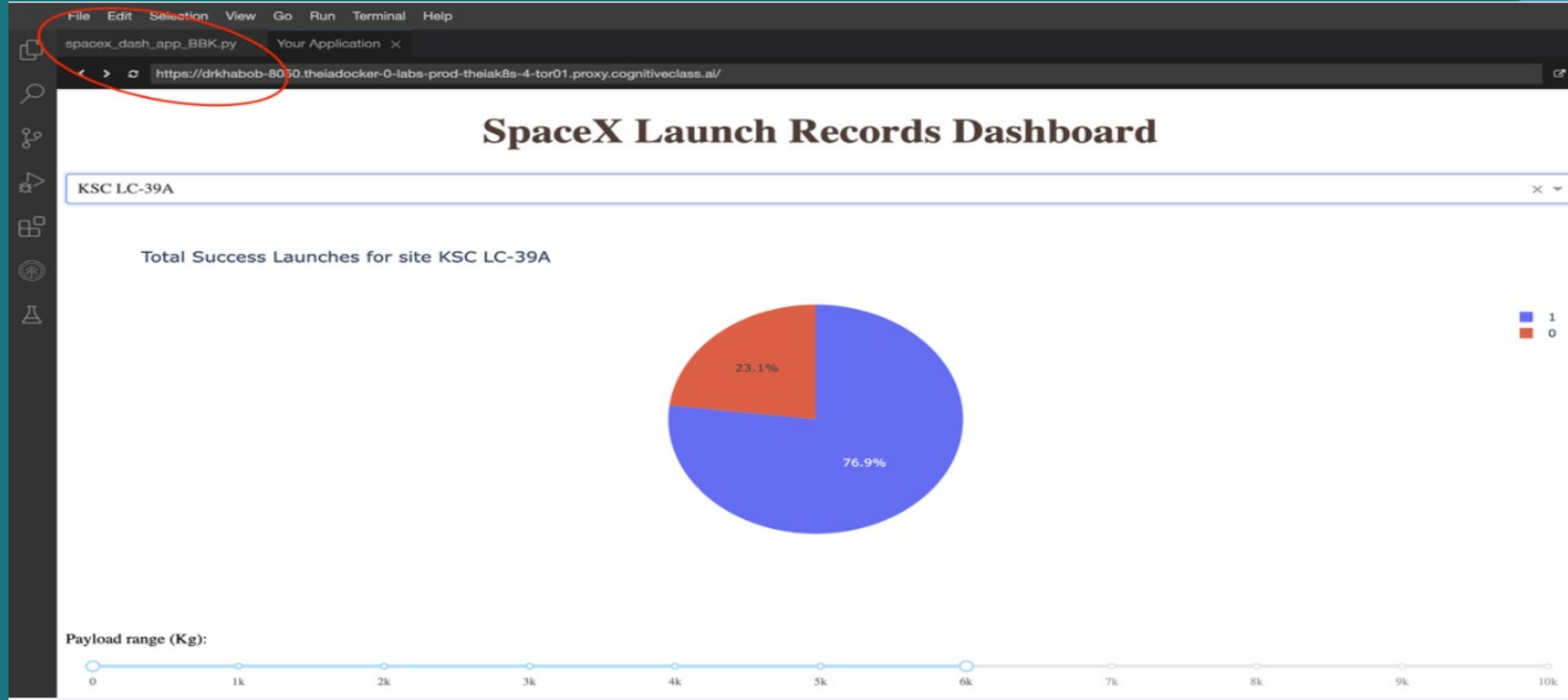
# Interactive analysis: Polygons to landmarks



# Dashboard: Launch success rate by launch site



# Dashboard: Success rate of the most prolific launch site

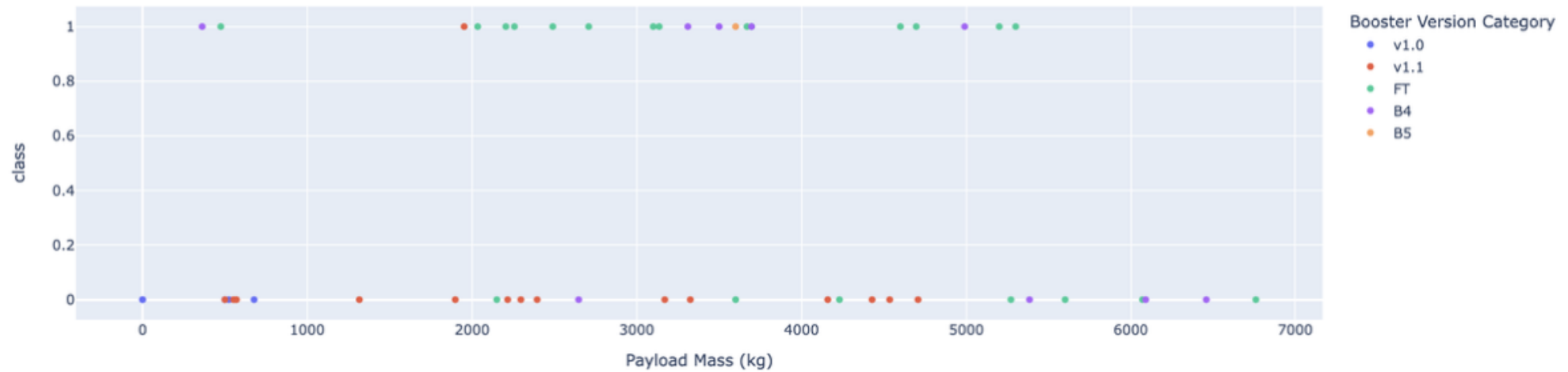


# Dashboard: Payload vs launch outcome

Payload range (Kg):



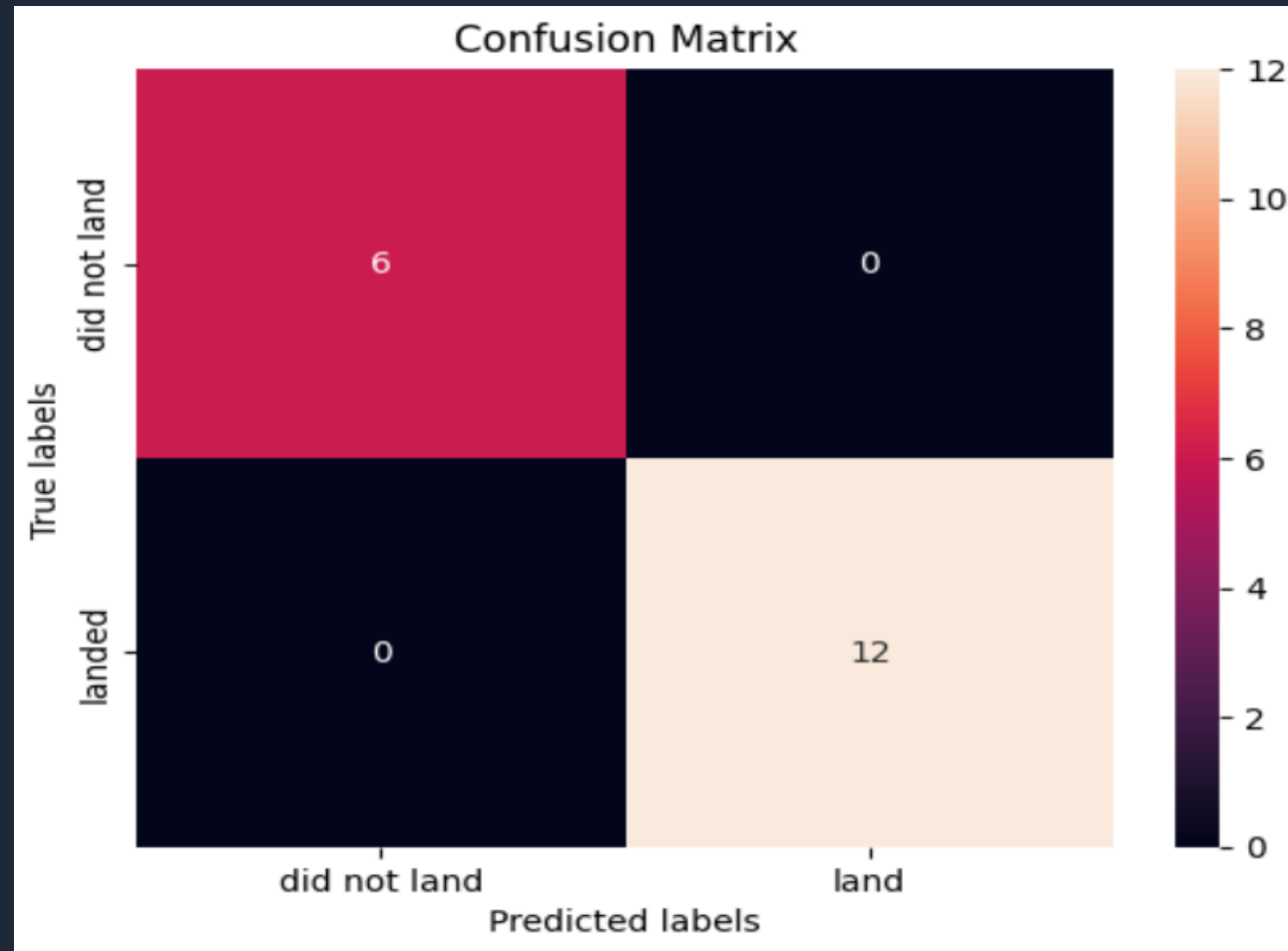
Launch Success Rate For All Sites



# Predictive analysis: classification accuracy



# Predictive analysis: K-nearest neighbor confusion matrix



# Conclusion

- Data science can be used to predict Falcon9 landings,
  - With near-perfect accuracy of out-of-sample prediction by some models i.e. KNN (and high accuracy by the other three models)
  - This can potentially save tens of millions in USD per launch
- Data mining of high-quality, clean and labelled data can be an asset when developing predictive models
- High-quality interactive dashboards and leaflet maps for day-to-day operations of business is demonstrably feasible
- Model evaluation methods to improve model selection based on meritocratic performance measures of out-of- sample data is also feasible





# Appendix

- Full github repository:  
<https://github.com/adityanaik240402/Applied-data-science-capstone-project>
- Falcon9 and Falcon Heavy Wikipedia page:  
[https://en.wikipedia.org/wiki/List\\_of\\_Falcon\\_9\\_and\\_Falcon\\_Heavy\\_launches](https://en.wikipedia.org/wiki/List_of_Falcon_9_and_Falcon_Heavy_launches)

