



IBM Developer  
SKILLS NETWORK

# Winning Space Race with Data Science

Srushti Anant Shimpi  
5th January 2025



# Outline

---

- Executive Summary
- Introduction
- Methodology
- Results
- Conclusion

# Executive Summary

---

- Summary of methodologies
  - Data Collection through API
  - Data Collection with Web Scraping
  - Data Wrangling
  - Exploratory Data Analysis with SQL
  - Exploratory Data Analysis with Data Visualization
  - Interactive Visual Analytics with Folium
  - Machine Learning Prediction
- Summary of all results
  - Exploratory Data Analysis result
  - Interactive analytics in screenshots
  - Predictive Analytics result

# 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. This goal of the project is to create a machine learning pipeline to predict if the first stage will land successfully.

- Problems you want to find answers

- What factors determine if the rocket will land successfully?
- The interaction amongst various features that determine the success rate of a successful landing.
- What operating conditions needs to be in place to ensure a successful landing program.



Section 1

# Methodology

# Methodology

---

## Executive Summary

- Data collection methodology:
  - Data was collected using SpaceX API and web scraping from Wikipedia.
- Perform data wrangling
  - One-hot encoding was applied to categorical features
- 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 model

# Data Collection

---

- The data was collected using various methods
  - Data collection was done using get request to the SpaceX API.
  - Next, we decoded the response content as a Json using `.json()` function call and turn it into a pandas dataframe using `.json_normalize()`.
  - We then cleaned the data, checked for missing values and fill in missing values where necessary.
  - In addition, we performed web scraping from Wikipedia for Falcon 9 launch records with BeautifulSoup.
  - The objective was to extract the launch records as HTML table, parse the table and convert it to a pandas dataframe for future analysis.



# Data Collection – SpaceX API

- We used the get request to the SpaceX API to collect data, clean the requested data and did some basic data wrangling and formatting.
- The link to the notebook is

<https://github.com/srushtishi/mpi/IBM-Data-Science-Capstone-SpaceX/blob/main/1 SpaceX Data Collection API.ipynb>

Now let's start requesting rocket launch data from SpaceX API with the following URL:

```
In [6]: spacex_url="https://api.spacexdata.com/v4/launches/past"
```

```
In [7]: response = requests.get(spacex_url)
```

Check the content of the response

```
In [8]: print(response.content)
```

```
b'{"fairings":{"reused":false,"recovery_attempt":false,"recovered":false,"s
es2.imgbox.com/94/f2/NN6Ph45r_o.png","large":{"https://images2.imgbox.com/5b/
unch":null,"media":null,"recovery":null},"flickr":{"small":[],"original":[]
om/watch?v=0a_00nJ_Y88","youtube_id":"0a_00nJ_Y88","article":{"https://www.sp
t-launch.html"},"wikipedia":{"https://en.wikipedia.org/wiki/DemoSat"},"static_
_fire_date_unix":1142553600,"net":false,"window":0,"rocket":"5e9d0d95eda6995
3,"altitude":null,"reason":"merlin engine failure"}],"details":"Engine failu
[],"ships":[],"capsules":[],"payloads":["5eb0e4b5b6c3bb0006eeb1e1"],"launchp
1,"name":"FalconSat","date_utc":"2006-03-24T22:30:00.000Z","date_unix":11432
ate_precision":"hour","upcoming":false,"cores":[{"core":"5e9e289df35918033d3
e,"reused":false,"landing_attempt":false,"landing_success":null,"landing_typ
d":false,"launch_library_id":null,"id":"5eb87cd9ffd86e000604b32a"},{"fairing
vered":false,"ships":[]},"links":{"patch":{"small":{"https://images2.imgbox.c
2.imgbox.com/80/a2/bkWotCIS_o.png"},"reddit":{"campaign":null,"launch":null,
[],"original":[]},"presskit":null,"webcast":{"https://www.youtube.com/watch?v
e":{"https://www.space.com/3590-spacex-falcon-1-rocket-fails-reach-orbit.html
```



# Data Collection - Scraping

- We applied web scrapping to webscrap Falcon 9 launch records with BeautifulSoup
- We parsed the table and converted it into a pandas dataframe.
- The link to the notebook is [https://github.com/srushtishimpi/IBM-Data-Science-Capstone-SpaceX/blob/main/2\\_spacex\\_Data\\_Collection\\_with\\_Web%20Scraping.ipynb](https://github.com/srushtishimpi/IBM-Data-Science-Capstone-SpaceX/blob/main/2_spacex_Data_Collection_with_Web%20Scraping.ipynb)

## TASK 1: Request the Falcon9 Launch Wiki page from its URL

First, let's perform an HTTP GET method to request the Falcon9 Launch HTML page, as an HTTP response.

```
[19]: # use requests.get() method with the provided static_url
      # assign the response to a object

      html_data = requests.get(static_url)
      html_data.status_code
```

```
it[19]: 200
```

Create a BeautifulSoup object from the HTML response

```
[20]: # Use BeautifulSoup() to create a BeautifulSoup object from a response text content

      soup = BeautifulSoup(html_data.text, "html.parser")
```

Print the page title to verify if the BeautifulSoup object was created properly

```
[21]: # Use soup.title attribute
      soup.title
```

```
it[21]: <title>List of Falcon 9 and Falcon Heavy launches - Wikipedia</title>
```

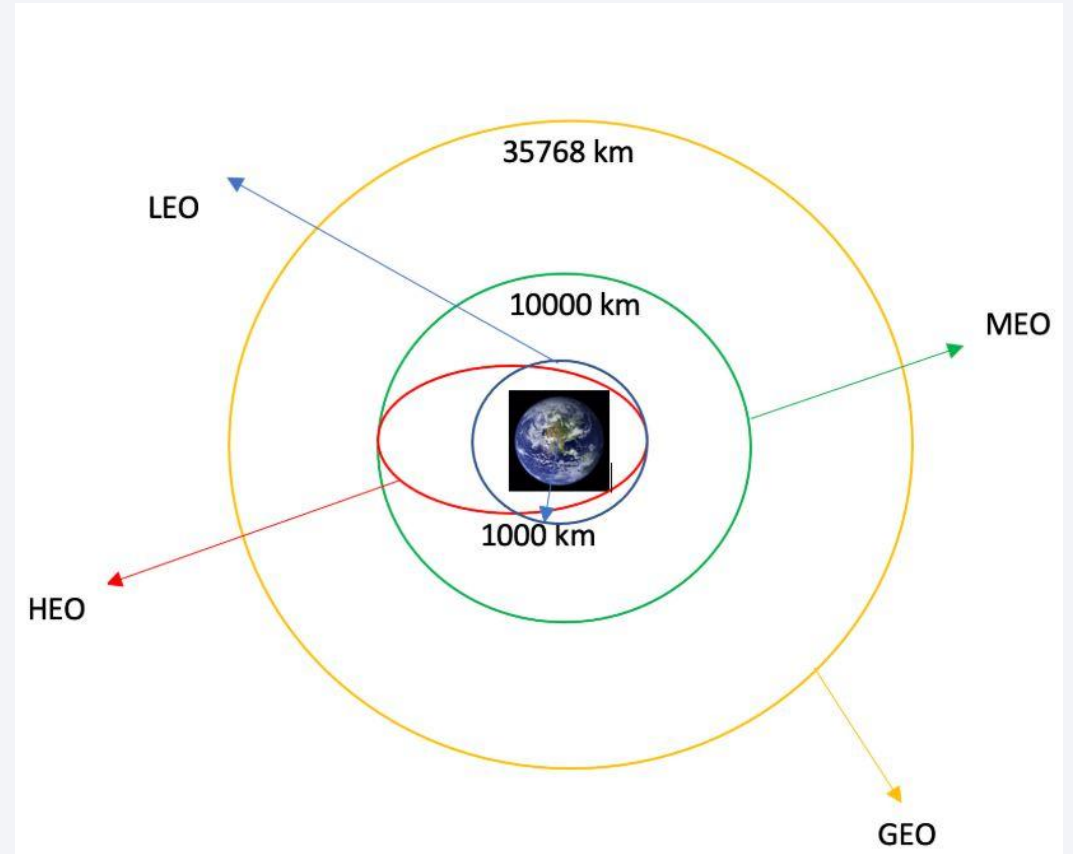
## TASK 2: Extract all column/variable names from the HTML table header

Next, we want to collect all relevant column names from the HTML table header

Let's try to find all tables on the wiki page first. If you need to refresh your memory about BeautifulSoup

# Data Wrangling

- We performed exploratory data analysis and determined the training labels.
- We calculated the number of launches at each site, and the number and occurrence of each orbits
- We created landing outcome label from outcome column and exported the results to csv.
- The link to the notebook is [https://github.com/srushtishimpi/IBM-Data-Science-Capstone-SpaceX/blob/main/3\\_spacex\\_Data\\_Wrangling.ipynb](https://github.com/srushtishimpi/IBM-Data-Science-Capstone-SpaceX/blob/main/3_spacex_Data_Wrangling.ipynb)



# EDA with Data Visualization

---

- We explored the data by visualizing the relationship between flight number and launch Site, payload and launch site, success rate of each orbit type, flight number and orbit type, the launch success yearly
- The link to the notebook is [https://github.com/srushtishimpi/IBM-Data-Science-Capstone-SpaceX/blob/main/5\\_SpaceX\\_EDA\\_with\\_Visualization.ipynb](https://github.com/srushtishimpi/IBM-Data-Science-Capstone-SpaceX/blob/main/5_SpaceX_EDA_with_Visualization.ipynb)

# EDA with SQL

---

- We loaded the SpaceX dataset into a PostgreSQL database without leaving the jupyter notebook.
- We applied EDA with SQL to get insight from the data. We wrote queries to find out for instance:
  - The names of unique launch sites in the space mission.
  - The total payload mass carried by boosters launched by NASA (CRS)
  - The average payload mass carried by booster version F9 v1.1
  - The total number of successful and failure mission outcomes
  - The failed landing outcomes in drone ship, their booster version and launch site names.
- The link to the notebook is [https://github.com/srushtishimpi/IBM-Data-Science-Capstone-SpaceX/blob/main/4\\_SpaceX\\_EDA\\_with\\_SQL.ipynb](https://github.com/srushtishimpi/IBM-Data-Science-Capstone-SpaceX/blob/main/4_SpaceX_EDA_with_SQL.ipynb)



# Build an Interactive Map with Folium

---

- We marked all launch sites, and added map objects such as markers, circles, lines to mark the success or failure of launches for each site on the folium map.
- We assigned the feature launch outcomes (failure or success) to class 0 and 1.i.e., 0 for failure, and 1 for success.
- Using the color-labeled marker clusters, we identified which launch sites have relatively high success rate.
- We calculated the distances between a launch site to its proximities. We answered some question for instance:
  - Are launch sites near railways, highways and coastlines.
  - Do launch sites keep certain distance away from cities.
- The link to the notebook is [https://github.com/srushtishimpi/IBM-Data-Science-Capstone-SpaceX/blob/main/6\\_SpaceX\\_Interactive\\_Visual\\_Analytics\\_with\\_Folium%20lab.ipynb](https://github.com/srushtishimpi/IBM-Data-Science-Capstone-SpaceX/blob/main/6_SpaceX_Interactive_Visual_Analytics_with_Folium%20lab.ipynb)

# Build a Dashboard with Plotly Dash

---

- We built an interactive dashboard with Plotly dash
- We plotted pie charts showing the total launches by a certain sites
- We plotted scatter graph showing the relationship with Outcome and Payload Mass (Kg) for the different booster version.
- The link to the notebook is [https://github.com/srushtishimpi/IBM-Data-Science-Capstone-SpaceX/blob/main/spacex\\_dash\\_app.py](https://github.com/srushtishimpi/IBM-Data-Science-Capstone-SpaceX/blob/main/spacex_dash_app.py)

# Predictive Analysis (Classification)

---

- We loaded the data using numpy and pandas, transformed the data, split our data into training and testing.
- We built different machine learning models and tune different hyperparameters using GridSearchCV.
- We used accuracy as the metric for our model, improved the model using feature engineering and algorithm tuning.
- We found the best performing classification model.
- The link to the notebook is  
[https://github.com/srushtishimpi/IBM-Data-Science-Capstone-SpaceX/blob/main/7\\_SpaceX\\_Machine%20Learning%20Prediction\\_Part\\_5.ipynb](https://github.com/srushtishimpi/IBM-Data-Science-Capstone-SpaceX/blob/main/7_SpaceX_Machine%20Learning%20Prediction_Part_5.ipynb)

# Results

---

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



The background of the slide is an abstract composition. It features a solid blue area on the left side, which transitions into a dynamic pattern of diagonal streaks in shades of blue and red on the right. These streaks are layered over a fine, light-colored grid, creating a sense of depth and movement, reminiscent of a digital or data visualization theme.

Section 2

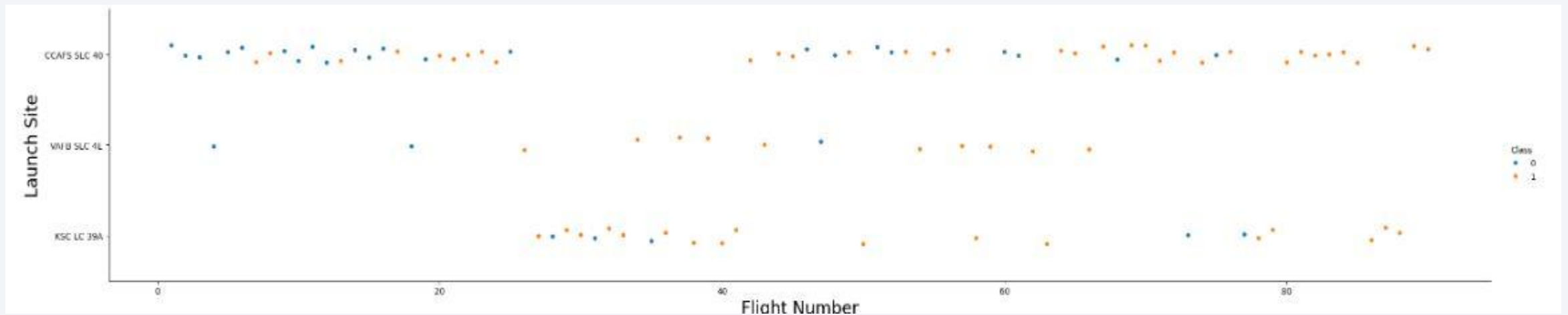
# Insights drawn from EDA



# Flight Number vs. Launch Site

---

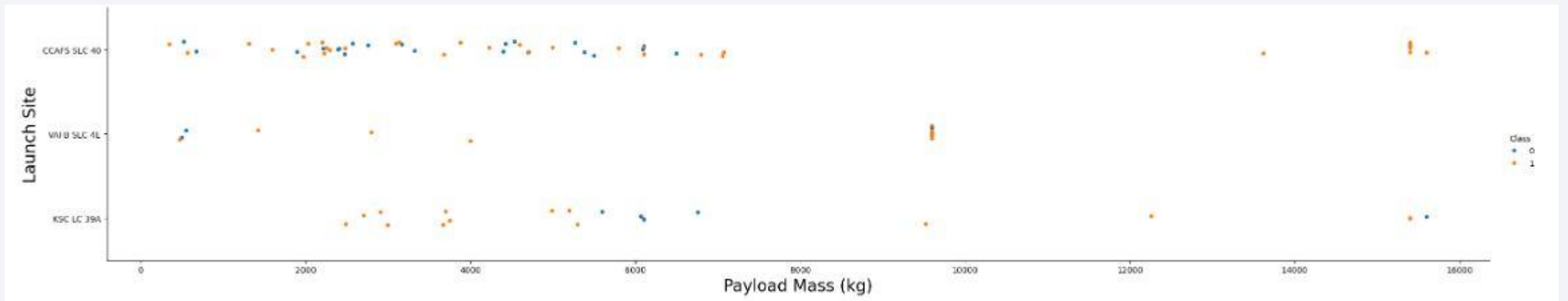
- From the plot, we found that the larger the flight amount at a launch site, the greater the success rate at a launch site.



# Payload vs. Launch Site

---

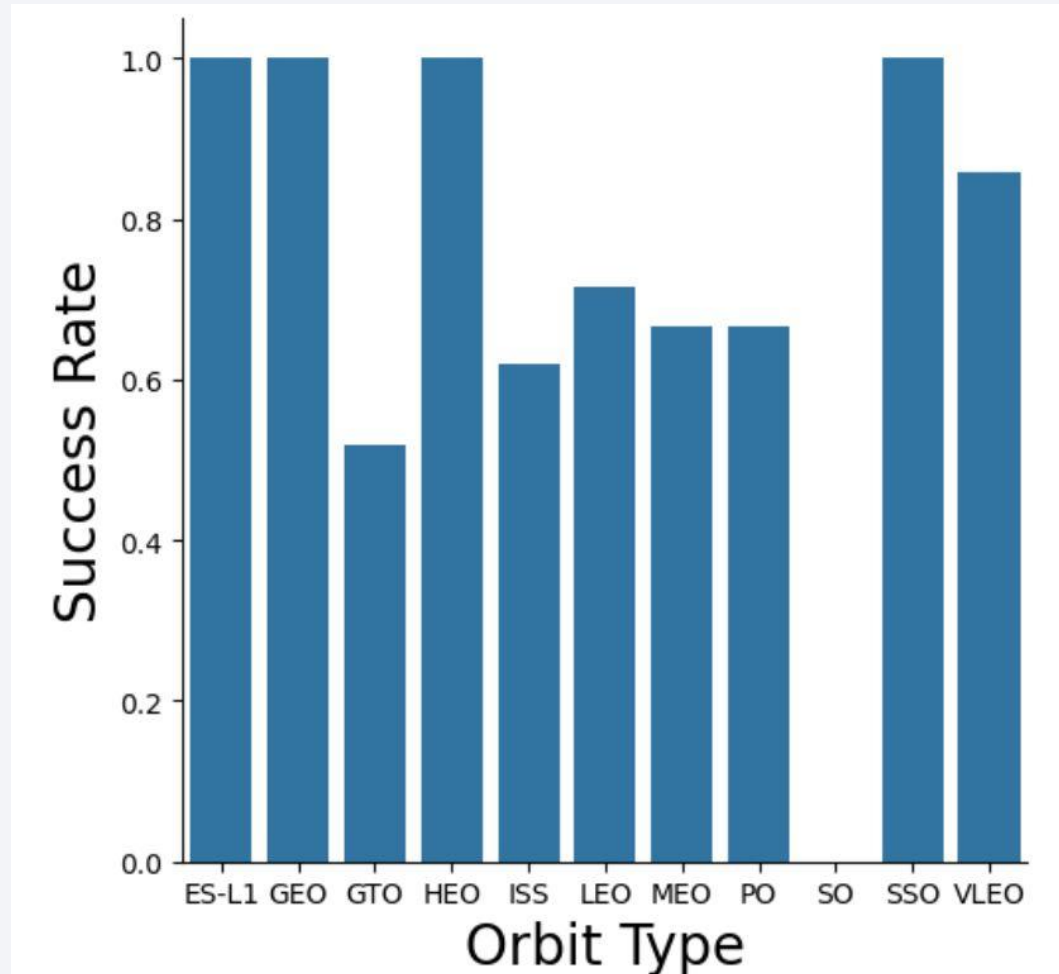
- The greater the Payload mass for launch site CCAFS SLC 40, higher the success rate for the rocket.



# Success Rate vs. Orbit Type

---

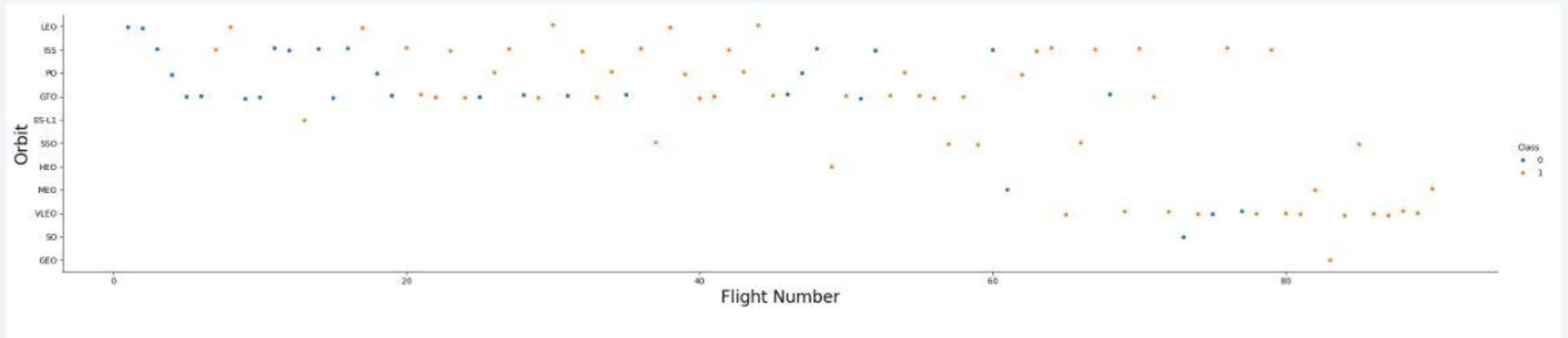
- From the plot, we can see that ES-L1, GEO, HEO, SSO had the most success rate.





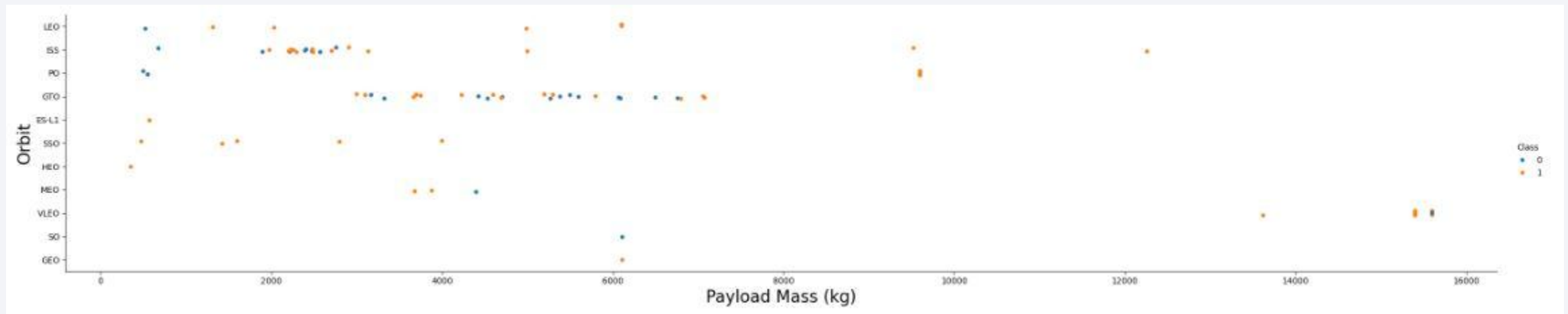
# Flight Number vs. Orbit Type

- The plot below shows the Flight Number vs. Orbit type. We observe that in the LEO orbit, success is related to the number of flights whereas in the GTO orbit, there is no relationship between flight number and the orbit.



# Payload vs. Orbit Type

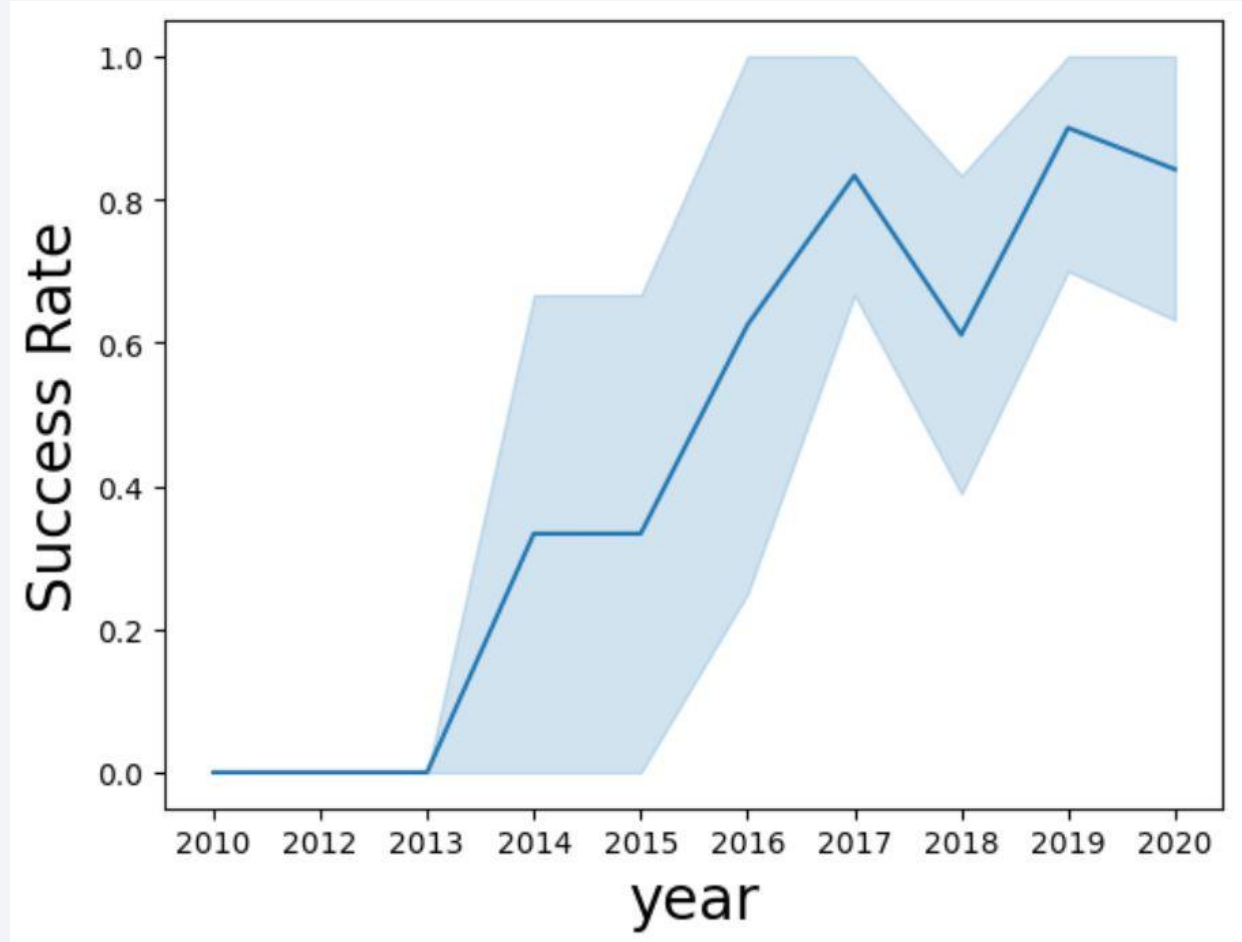
- We can observe that with heavy payloads, the successful landing are more for PO, LEO and ISS orbits.



# Launch Success Yearly Trend

---

- From the plot, we can observe that success rate since 2013 kept on increasing till 2020.



# All Launch Site Names

---

- We used the keyword **DISTINCT** to show only unique launch sites from the SpaceX data.

## Task 1

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

```
%sql select distinct launch_site from SPACEXTABLE;
```

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

### Launch\_Site

CCAFS LC-40

VAFB SLC-4E

KSC LC-39A

CCAFS SLC-40



# Launch Site Names Begin with 'CCA'

- We used the this query to display 5 records where launch sites begin with 'CCA'

## Task 2

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

```
%sql select * from SPACEXTABLE 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

---

- We calculated the total payload carried by boosters from NASA as 45596 using the query below

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

```
%sql select sum(payload_mass__kg_) as total_payload_mass from SPACEXTABLE where customer = 'NASA (CRS)';
```

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

```
Done.
```

<b>total_payload_mass</b>
---------------------------

45596
-------

# Average Payload Mass by F9 v1.1

---

- We calculated the average payload mass carried by booster version F9 v1.1 as 2534.66

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

```
%sql select avg(payload_mass__kg_) as average_payload_mass from SPACEXTABLE where booster_version like '%F9 v1.1%';
```

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

```
Done.
```

<u>average_payload_mass</u>
-----------------------------

2534.6666666666665
--------------------

# First Successful Ground Landing Date

---

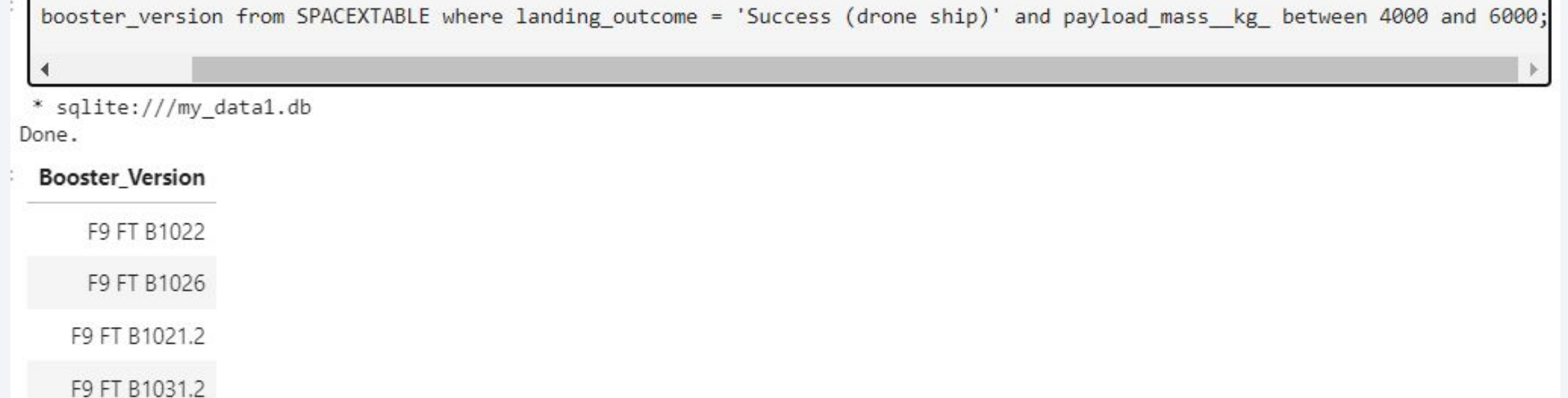
- We observed that the dates of the first successful landing outcome on ground pad was 22<sup>nd</sup> December 2015

```
%sql select min(date) as first_successful_landing from SPACEXTABLE where landing_outcome = 'Success (ground pad)';  
* sqlite:///my_data1.db  
Done.  
  
first_successful_landing  
2015-12-22
```

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

- We used the **WHERE** clause to filter for boosters which have successfully landed on drone ship and applied the **AND** condition to determine successful landing with payload mass greater than 4000 but less than 6000

```
: booster_version from SPACEXTABLE where landing_outcome = 'Success (drone ship)' and payload_mass__kg_ between 4000 and 6000;
```



The screenshot shows a SQLite terminal window. At the top, a SQL query is entered: `booster_version from SPACEXTABLE where landing_outcome = 'Success (drone ship)' and payload_mass__kg_ between 4000 and 6000;`. Below the query bar, the terminal shows the connection path `* sqlite:///my_data1.db` and the status `Done.`. The results are displayed as a table with a single column header **Booster\_Version**. The table contains four rows of data: `F9 FT B1022`, `F9 FT B1026`, `F9 FT B1021.2`, and `F9 FT B1031.2`. The second and fourth rows are highlighted with a light gray background.

Booster_Version
F9 FT B1022
F9 FT B1026
F9 FT B1021.2
F9 FT B1031.2



# Total Number of Successful and Failure Mission Outcomes

- We used wildcard like '%' to filter for **WHERE** Mission Outcome was a success or a failure.

```
%sql select mission_outcome, count(*) as total_number from SPACEXTABLE group by mission_outcome;
```

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

```
Done.
```

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

# Boosters Carried Maximum Payload

```
%sql select booster_version from SPACESTABLE where payload_mass_kg_ = (select max(payload_mass_kg_) from SPACESTABLE);
```

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

Booster_Version
-----------------

F9 B5 B1048.4
---------------

F9 B5 B1049.4
---------------

F9 B5 B1051.3
---------------

F9 B5 B1056.4
---------------

F9 B5 B1048.5
---------------

F9 B5 B1051.4
---------------

F9 B5 B1049.5
---------------

F9 B5 B1060.2
---------------

F9 B5 B1058.3
---------------

F9 B5 B1051.6
---------------

F9 B5 B1060.3
---------------

F9 B5 B1049.7
---------------

- We determined the booster that have carried the maximum payload using a subquery in the **WHERE** clause and the **MAX()** function.

# 2015 Launch Records

---

- We used a combinations of the **WHERE** clause, **LIKE**, **AND**, and **BETWEEN** conditions to filter for failed landing outcomes in drone ship, their booster versions, and launch site names for year 2015

```
%%sql select booster_version, launch_site, landing_outcome
from SPACEXTABLE
where landing_outcome = 'Failure (drone ship)'
AND Date BETWEEN '2015-01-01' AND '2015-12-31';
```

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

Booster_Version	Launch_Site	Landing_Outcome
F9 v1.1 B1012	CCAFS LC-40	Failure (drone ship)
F9 v1.1 B1015	CCAFS LC-40	Failure (drone ship)

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

```
%%sql select landing_outcome, count(*) as count_outcomes from SPACEXTABLE
      where date between '2010-06-04' and '2017-03-20'
      group by landing_outcome
      order by count_outcomes desc;
```

\* sqlite:///my\_data1.db  
Done.

Landing_Outcome	count_outcomes
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

- We selected Landing outcomes and the **COUNT** of landing outcomes from the data and used the **WHERE** clause to filter for landing outcomes **BETWEEN** 2010-06-04 to 2017-03-20.
- We applied the **GROUP BY** clause to group the landing outcomes and the **ORDER BY** clause to order the grouped landing outcome in descending order.

A satellite view of Earth from space, showing the curvature of the planet and city lights at night. The image is a composite of a dark blue sky with stars and a view of the Earth's surface from space. The Earth's surface is mostly dark blue, with a thin layer of white clouds. A bright, glowing arc of city lights is visible along the horizon, indicating a coastal or urban area. The text "Section 3" is overlaid on the left side of the image.

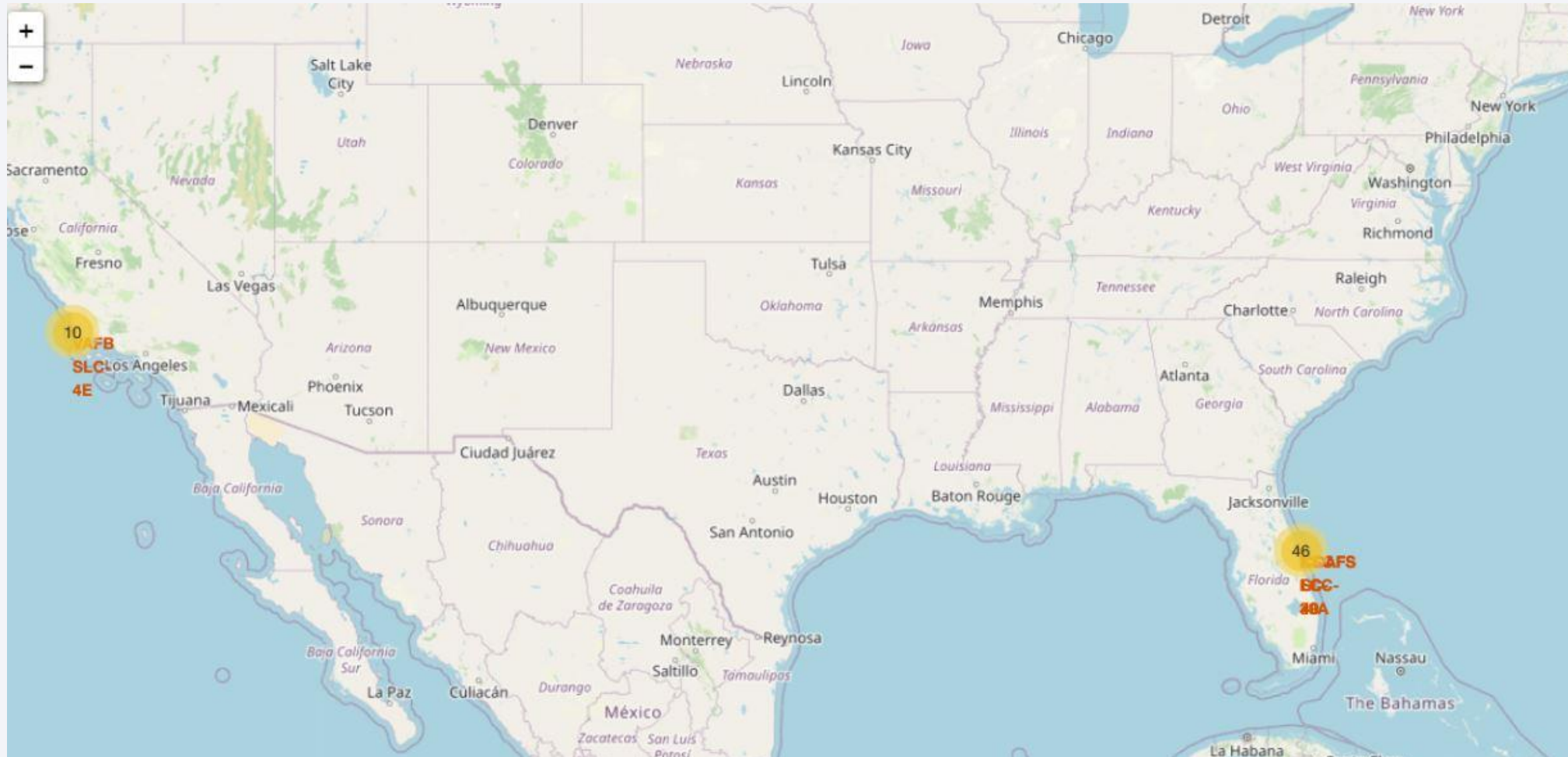
Section 3

# Launch Sites Proximities Analysis



# All launch sites global map markers

This map markers shows NASA launch sites in California and Florida state.



# Markers showing launch sites with color labels

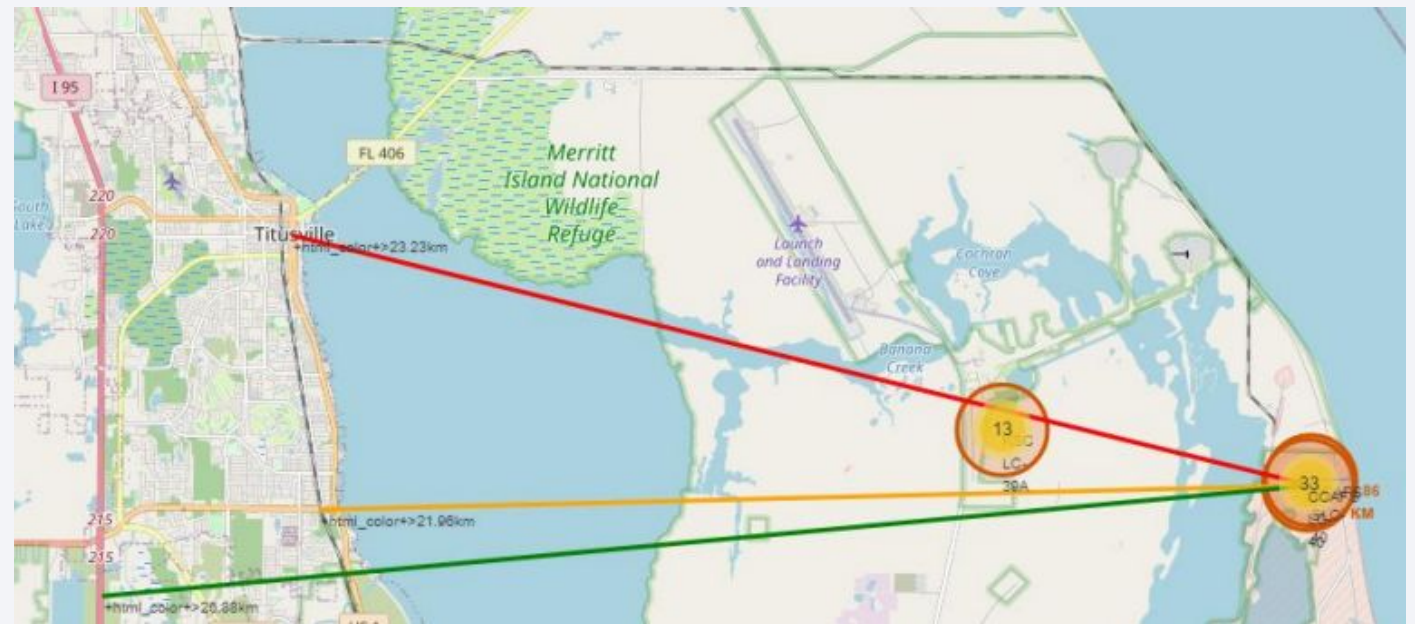


- Green markers for successful launches
- Red markers for unsuccessful launches
- Launch site **CCAFS SLC-40** has a **3/7 success rate (42.9%)**

# Distance to Proximities

- **CCAFS SLC-40**

- 0.86 km from nearest coastline
- 21.96 km from nearest railway
- 23.23 km from nearest city
- 26.88 km from nearest highway







Section 4

# Build a Dashboard with Plotly Dash

# Launch Success by Site

---

- KSC LC-39A has the most successful launches amongst all launch sites (41.2 %)

Total Success Launches by Site



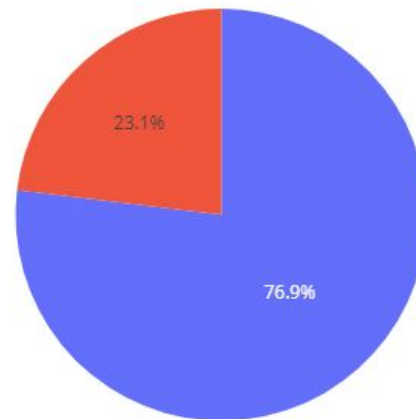


# Launch Success (KSC LC-39A)

---

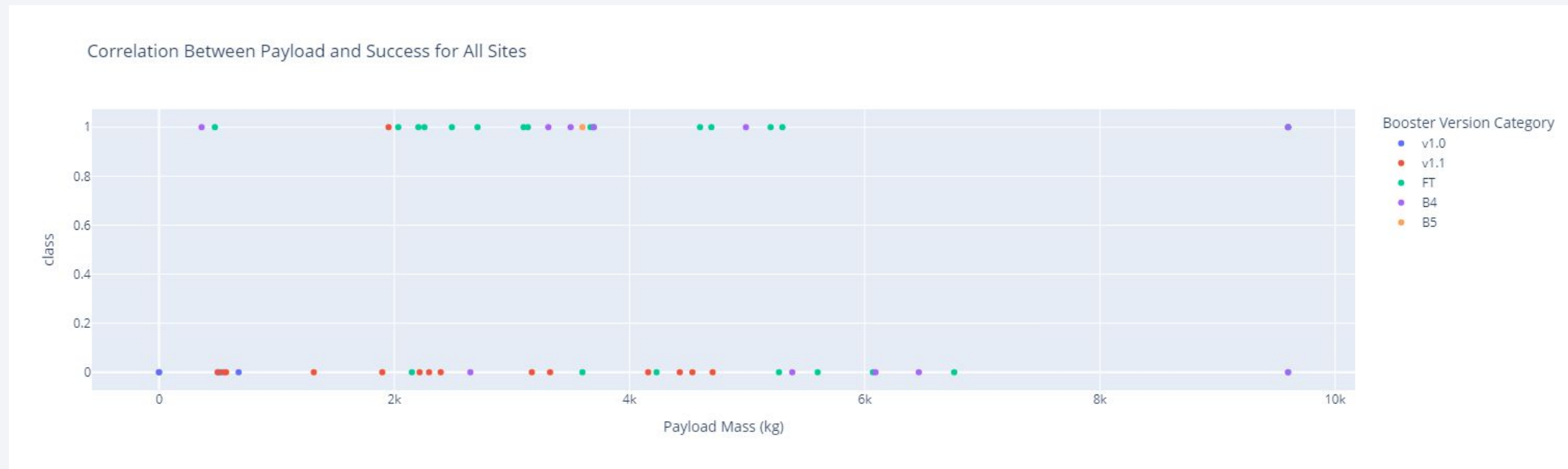
- KSC LC-39A has the highest success rate amongst launch sites (76.9%)
- 10 successful launches and 3 failed launches

Total Success Launches for Site KSC LC-39A



# Payload VS Launch Success Outcome for all sites

- Payloads between 2,000 kg and 5,000 kg have the highest success rate
- 1 indicating successful outcome and 0 indicating an unsuccessful outcome





Section 5

# Predictive Analysis (Classification)

# Classification Accuracy

---

- The decision tree classifier is the model with the highest classification accuracy (87.85%)

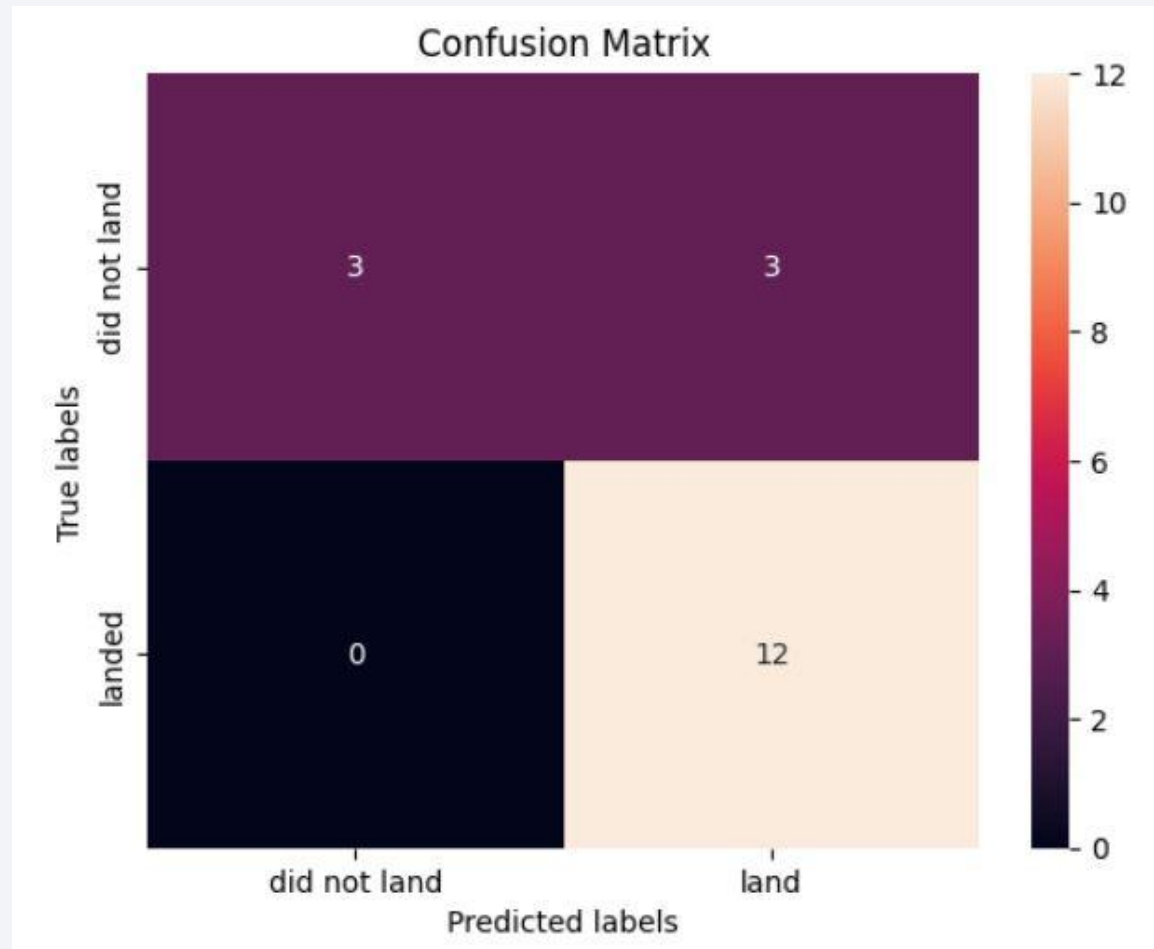
```
models = {'KNeighbors': knn_cv.best_score_,
          'DecisionTree': tree_cv.best_score_,
          'LogisticRegression': logreg_cv.best_score_,
          'SupportVector': svm_cv.best_score_}

bestalgorithm = max(models, key=models.get)
print('Best model is', bestalgorithm, 'with a score of', models[bestalgorithm])
if bestalgorithm == 'DecisionTree':
    print('Best params is :', tree_cv.best_params_)
if bestalgorithm == 'KNeighbors':
    print('Best params is :', knn_cv.best_params_)
if bestalgorithm == 'LogisticRegression':
    print('Best params is :', logreg_cv.best_params_)
if bestalgorithm == 'SupportVector':
    print('Best params is :', svm_cv.best_params_)
```

Best model is DecisionTree with a score of 0.8785714285714284

Best params is : {'criterion': 'entropy', 'max\_depth': 4, 'max\_features': 'sqrt', 'min\_samples\_leaf': 1, 'min\_samples\_split': 2, 'splitter': 'random'}

# Confusion Matrix



- The confusion matrix for the decision tree classifier shows that the classifier can distinguish between the different classes. The major problem is the false positives .i.e., unsuccessful landing marked as successful landing by the classifier.

# Conclusions

---

We can conclude that:

- The larger the flight amount at a launch site, the greater the success rate at a launch site.
- Launch success rate started to increase in 2013 till 2020.
- Orbits ES-L1, GEO, HEO, SSO, VLEO had the most success rate.
- KSC LC-39A had the most successful launches of any sites.
- The Decision tree classifier is the best machine learning algorithm for this task.



Thank you!

