



IBM Developer  
SKILLS NETWORK

# Winning Space Race with Data Science

Mani Kumar Raju Kovvuru  
25/04/2022



A large orange circle is positioned on the left side of the slide, partially cut off by the edge. The word "Outline" is written in white text inside this circle.

# Outline

---

Executive Summary

---

Introduction

---

Methodology

---

Results

---

Conclusion

---

Appendix

# Executive Summary

## Summary of methodologies

- Data Collection with API
- Data Collection by Web Scraping
- Data Wrangling
- Exploratory Data Analysis with SQL
- Exploratory Data Analysis with Data Visualization
- Interactive Data Visualization with Folium
- Machine Learning Predict

## Summary of all results

- Exploratory Data Analysis
- Analytics in screenshots
- Predictive Analytics with Machine Learning result

# Introduction

## Project background and context

- Space X 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 Space X 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 space X 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

Click to add text

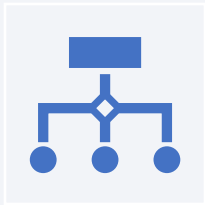
# Methodology

- Executive Summary
- Data collection methodology:
  - Data is collected using SpaceX API and web scraping from Wikipedia.
- Perform data wrangling
  - One-hot encoding is 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 models

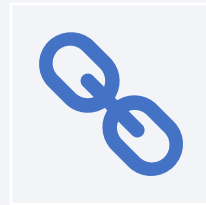
# Data Collection

- 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/manikumar261996/Capstone\\_RocketLaunch\\_Analysis/blob/master/DataCollection%20API.ipynb](https://github.com/manikumar261996/Capstone_RocketLaunch_Analysis/blob/master/DataCollection%20API.ipynb)

```
1. Get request for rocket launch data using API

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

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

2. Use json_normalize method to convert json result to dataframe

In [12]: # Use json_normalize method to convert the json result into a dataframe
         # decode response content as json
         static_json_df = res.json()

In [13]: # apply json_normalize
         data = pd.json_normalize(static_json_df)

3. We then performed data cleaning and filling in the missing values

In [30]: rows = data_falcon9['PayloadMass'].values.tolist()[0]

         df_rows = pd.DataFrame(rows)
         df_rows = df_rows.replace(np.nan, PayloadMass)

         data_falcon9['PayloadMass'][0] = df_rows.values
         data_falcon9
```



# 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/manikumar261996/Capstone\\_RocketLaunch\\_Analyis/blob/master/WebScrapping.ipynb](https://github.com/manikumar261996/Capstone_RocketLaunch_Analyis/blob/master/WebScrapping.ipynb)

```
1. Apply HTTP Get method to request the Falcon 9 rocket launch page

In [4]: static_url = "https://en.wikipedia.org/w/index.php?title=List_of_Falcon_9_and_Falcon_Heavy_launches&oldid=1027686922"

In [5]: # 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

Out[5]: 200

2. Create a BeautifulSoup object from the HTML response

In [6]: # 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

In [7]: # Use soup.title attribute
soup.title

Out[7]: <title>List of Falcon 9 and Falcon Heavy launches - Wikipedia</title>

3. Extract all column names from the HTML table header

In [10]: column_names = []

# Apply find_all() function with 'th' element on first_launch_table
# Iterate each th element and apply the provided extract_column_from_header() to get a column name
# Append the Non-empty column name ('if name is not None and len(name) > 0') into a list called column_names

element = soup.find_all('th')
for row in range(len(element)):
    try:
        name = extract_column_from_header(element[row])
        if (name is not None and len(name) > 0):
            column_names.append(name)
    except:
        pass

4. Create a dataframe by parsing the launch HTML tables
5. Export data to csv
```

# 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/manikumar261996/Capstone\\_RocketLaunch\\_Analysis/blob/master/EDA\\_Data\\_Wrangling.ipynb](https://github.com/manikumar261996/Capstone_RocketLaunch_Analysis/blob/master/EDA_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 trend.
- The link to the notebook is [https://github.com/manikumar261996/Capstone\\_RocketLaunch\\_Analyis/blob/master/EDA\\_Visualization.ipynb](https://github.com/manikumar261996/Capstone_RocketLaunch_Analyis/blob/master/EDA_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/manikumar261996/Capstone\\_RocketLaunch\\_Analysis/blob/master/EDA\\_SQL.ipynb](https://github.com/manikumar261996/Capstone_RocketLaunch_Analysis/blob/master/EDA_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.



# 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/manikumar261996/Capstone\\_RocketLaunch\\_Analyis/blob/master/spacex\\_dash\\_app.py](https://github.com/manikumar261996/Capstone_RocketLaunch_Analyis/blob/master/spacex_dash_app.py)
- [https://github.com/manikumar261996/Capstone\\_RocketLaunch\\_Analyis/blob/master/VisualAnalytics\\_FoliumLab.ipynb](https://github.com/manikumar261996/Capstone_RocketLaunch_Analyis/blob/master/VisualAnalytics_FoliumLab.ipynb)

# 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/manikumar261996/Capstone\\_RocketLaunch\\_Analyis/blob/master/ML\\_Prediction\\_LaunchAnalysis.ipynb](https://github.com/manikumar261996/Capstone_RocketLaunch_Analyis/blob/master/ML_Prediction_LaunchAnalysis.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 dark blue field on the left side, which transitions into a complex pattern of diagonal streaks in shades of blue, red, and teal on the right. These streaks have a textured, almost woven appearance. Overlaid on this pattern is a faint, light blue grid that recedes into the distance, creating a sense of depth and perspective.

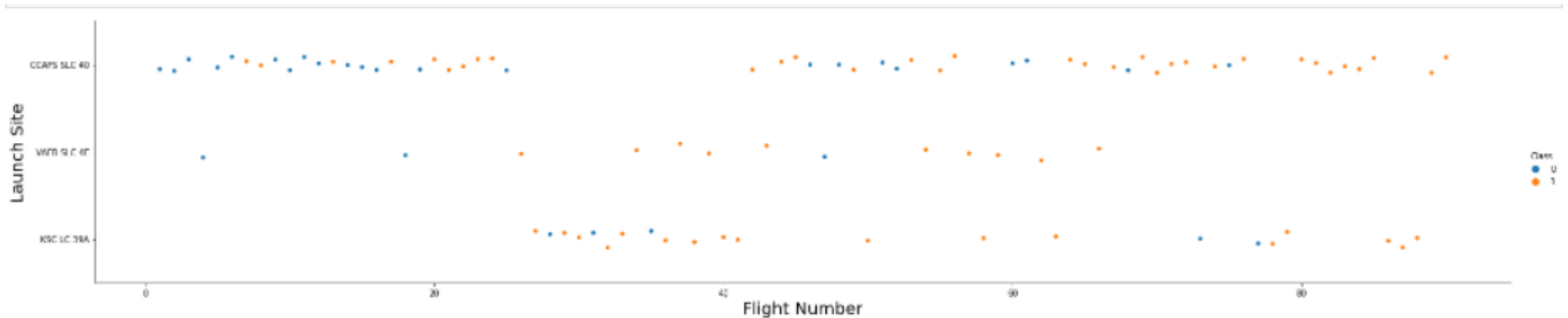
Section 2

# Insights drawn from EDA



## Flight Number vs. Launch Site

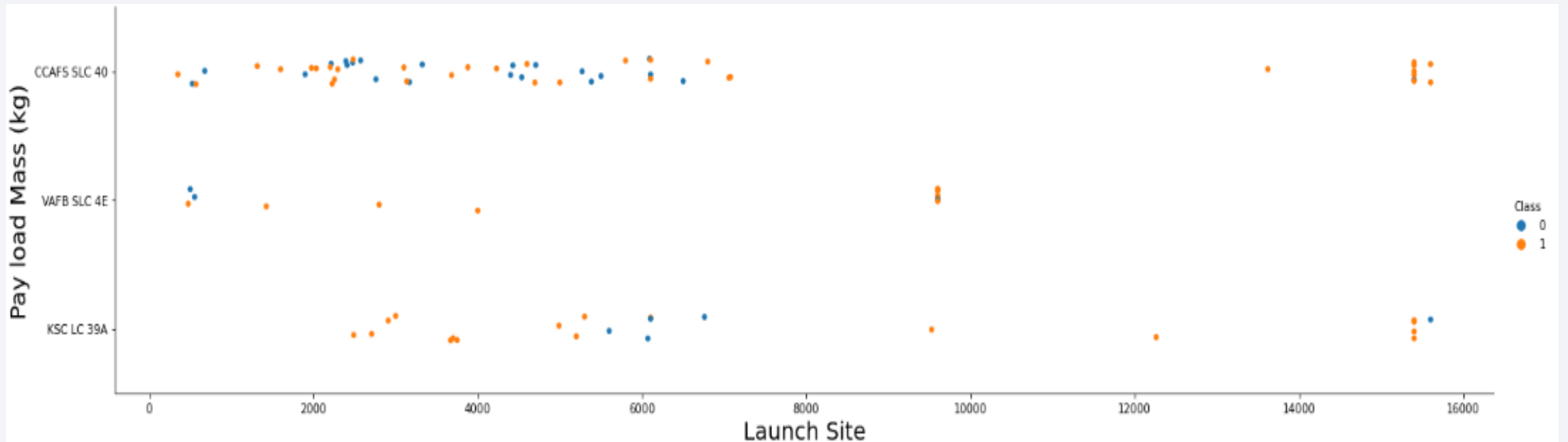
From the plot, we found that the larger the flight amount at a launch site, the greater the success rate at a launch site.





# Payload vs. Launch Site

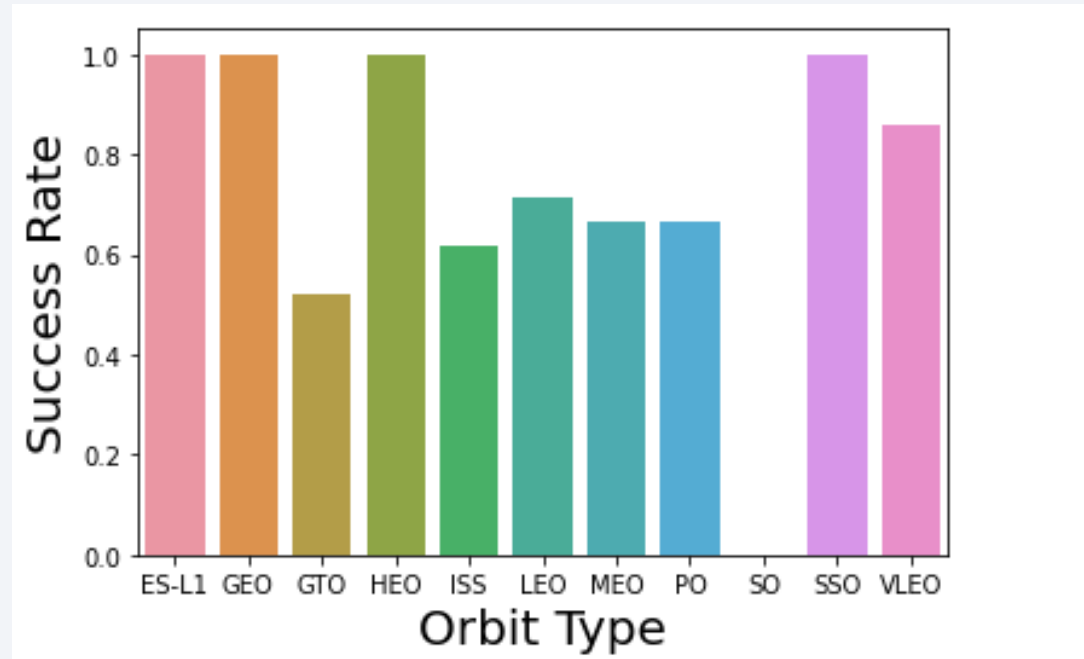
- The greater the payload for CCAFS SLC 40 the higher the success rate for Rocket Launch



# Success Rate vs. Orbit Type

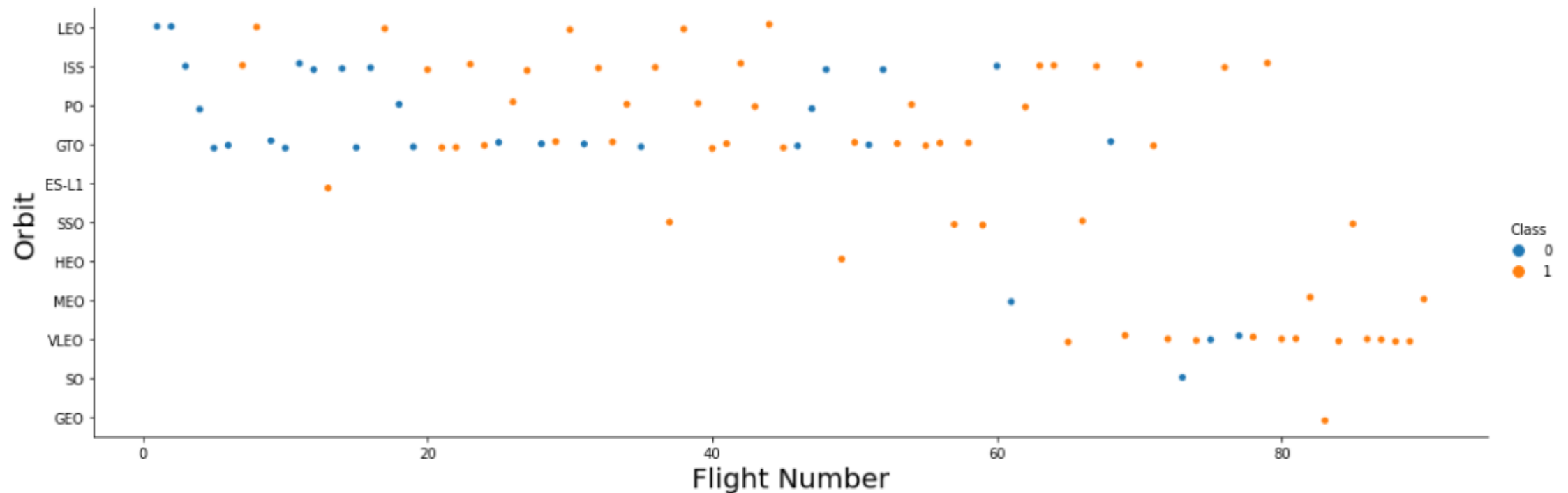
---

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



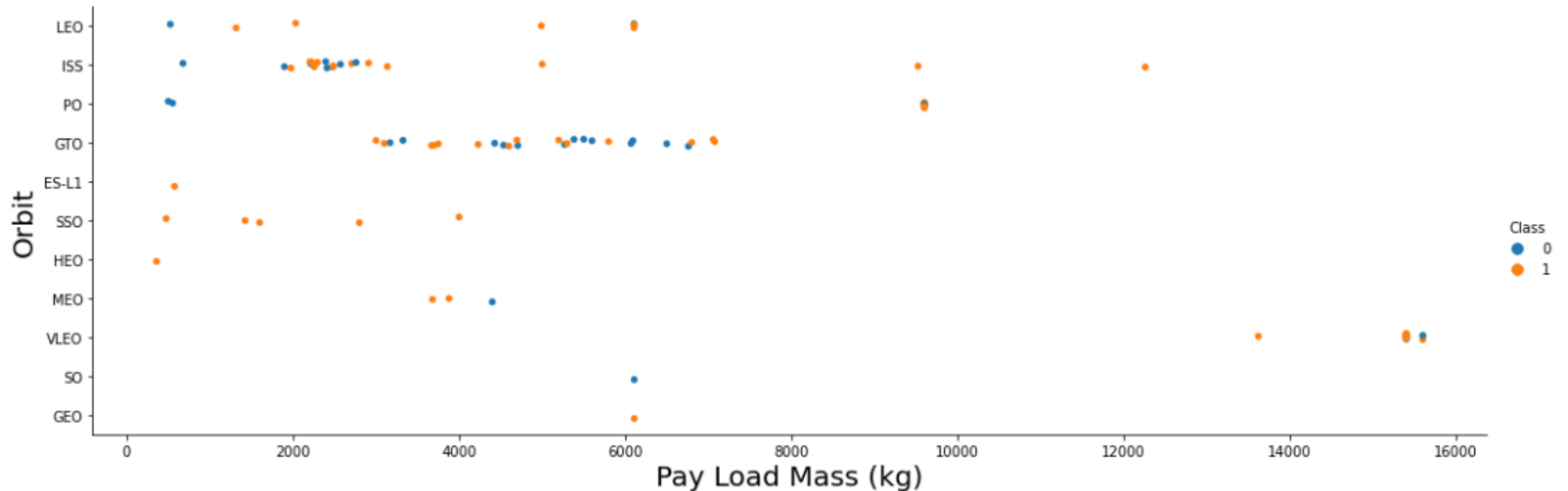
## Flight Number vs. Orbit Type

- 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

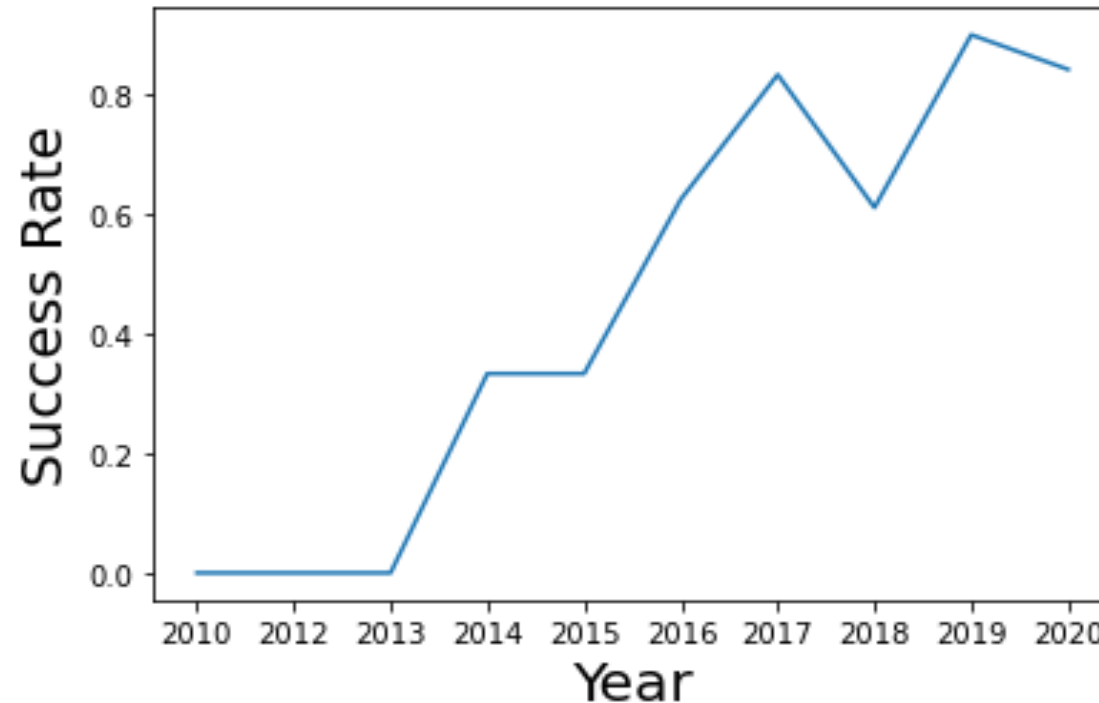
- We can observe that with heavy payloads, the successful landings are more for PO, LEO and ISS orbits.



# Launch Success Yearly Trend

---

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



you can observe that the success rate since 2013 kept increasing till 2020



# All Launch Site Names

---

- We used the key word **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

```
4]: %sql select distinct LAUNCH_SITE from SPACEXTBL
* ibm_db_sa://vck08231:***@55fbc997-9266-4331-afd3-888b05e734c0.bs2io90l08kqb1od8lcg.databases.appdomain.cloud:31929/bludb
Done.
4]: launch_site
CCAFS LC-40
CCAFS SLC-40
KSC LC-39A
VAFB SLC-4E
```

# Launch Site Names Begin with 'CCA'

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

## Task 2

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

In [7]: `%sql select * from SPACEXTBL where LAUNCH_SITE like 'CCA%' LIMIT 5`

`* ibm_db_sa://vck08231:***@55fbc997-9266-4331-afd3-888b05e734c0.bs2io90108kqb1od81cg.databases.appdomain.cloud:31929/bludb`  
Done.

Out[7]:

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

---

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

## Task 3

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

```
In [23]: %sql select sum(payload_mass__kg_) as Total_Payload, customer from SPACEXTBL where customer = 'NASA (CRS)' group by customer
```

```
* ibm_db_sa://vck08231:***@55fbc997-9266-4331-afd3-888b05e734c0.bs2io90108kqb1od8lcg.databases.appdomain.cloud:31929/bludb
Done.
```

```
Out[23]: total_payload  customer
```

```
45596  NASA (CRS)
```

# Average Payload Mass by F9 v1.1

---

- We calculated the average payload mass carried by booster version F9 v1.1 as 2928.4

## Task 4

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

```
In [28]: %sql select avg(payload_mass_kg_) as Avg_Payload, booster_version from SPACEXTBL where booster_version like 'F9 v1.1' group by booster_version
```

```
* ibm_db_sa://vck08231:***@55fbc997-9266-4331-afd3-888b05e734c0.bs2io90l08kqb1od8lcg.databases.appdomain.cloud:31929/bludb  
Done.
```

```
Out[28]: avg_payload  booster_version
```

2928	F9 v1.1
------	---------

# First Successful Ground Landing Date

---

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

## Task 5

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

*Hint: Use min function*

```
] : %sql select min(DATE) as First_Success from SPACEXTBL where LANDING__OUTCOME='Success (ground pad)'  
  
* ibm_db_sa://vck08231:***@55fbc997-9266-4331-afd3-888b05e734c0.bs2io90l08kqb1od8lcg.databases.appdomain.cloud:31929/bludb  
Done.  
]: first_success  
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

### Task 6

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

```
In [35]: %sql select BOOSTER_VERSION from SPACEXTBL where PAYLOAD_MASS__KG_ between 4000 and 6000 AND LANDING__OUTCOME='Success (drone ship)'
```

\* ibm\_db\_sa://vck08231:\*\*\*@55fbc997-9266-4331-afd3-888b05e734c0.bs2io90l08kqb1od8lcg.databases.appdomain.cloud:31929/bludb  
Done.

```
Out[35]: 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** MissionOutcome was a success or a failure.

## Task 7

List the total number of successful and failure mission outcomes

```
In [37]: %sql SELECT COUNT(*) as Total_Outcomes from SPACEXTBL where MISSION_OUTCOME like '%Success%' or MISSION_OUTCOME like '%Failure%'

* ibm_db_sa://vck08231:***@55fbc997-9266-4331-afd3-888b05e734c0.bs2io90l08kqb1od8l1cg.databases.appdomain.cloud:31929/bludb
Done.
```

```
Out[37]: total_outcomes
```

```
101
```

# Boosters Carried Maximum Payload

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

## Task 8

List the names of the booster\_versions which have carried the maximum payload mass. Use a subquery

```
In [39]: %sql SELECT BOOSTER_VERSION,PAYLOAD_MASS_KG_ from SPACEXTBL where PAYLOAD_MASS_KG_ = (SELECT MAX(PAYLOAD_MASS_KG_) from SPACEXTBL)
```

\* ibm\_db\_sa://vck08231:\*\*\*@55fbc997-9266-4331-afd3-888b05e734c0.bs2io90l08kqb1od8lcg.databases.appdomain.cloud:31929/bludb  
Done.

```
Out[39]:
```

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 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

## Task 9

List the failed landing\_outcomes in drone ship, their booster versions, and launch site names for in year 2015

```
[44]: %sql select landing__outcome,booster_version,launch_site,DATE from SPACEXTBL where landing__outcome = 'Failure (drone ship)' and year(DATE)=2015

* ibm_db_sa://vck08231:***@55fbc997-9266-4331-afd3-888b05e734c0.bs2io90l08kqb1od8lcg.databases.appdomain.cloud:31929/bludb
Done.
```

landing_outcome	booster_version	launch_site	DATE
Failure (drone ship)	F9 v1.1 B1012	CCAFS LC-40	2015-01-10
Failure (drone ship)	F9 v1.1 B1015	CCAFS LC-40	2015-04-14

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

- 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.

## Task 10

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

```
In [49]: %sql select count(landing__outcome) as count ,landing__outcome from SPACEXTBL where Date between '2010-06-04' and '2017-03-20' group by landing__outcom

* ibm_db_sa://vck08231:***@55fbc997-9266-4331-afd3-888b05e734c0.bs2io90l08kqb1od8l1cg.databases.appdomain.cloud:31929/bludb
Done.
```

```
Out[49]:
```

COUNT	landing__outcome
10	No attempt
5	Failure (drone ship)
5	Success (drone ship)
3	Controlled (ocean)
3	Success (ground pad)
2	Failure (parachute)
2	Uncontrolled (ocean)
1	Precluded (drone ship)

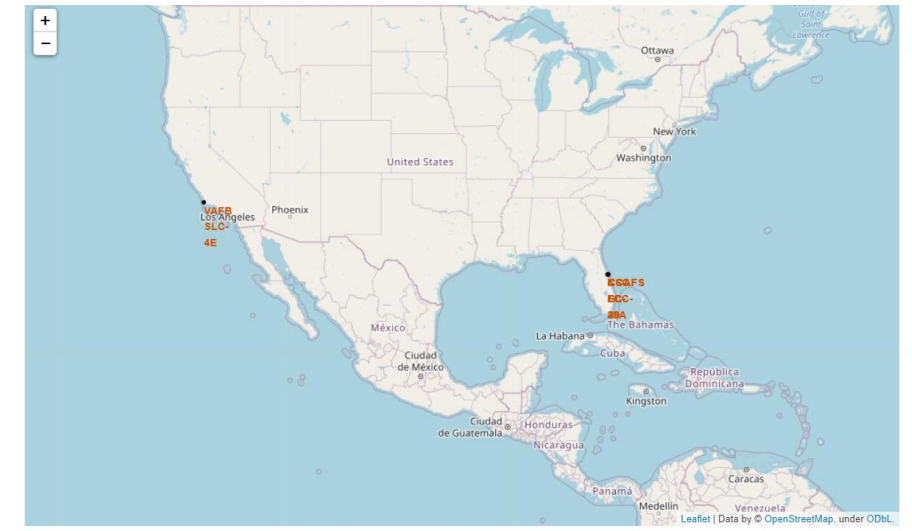
A satellite view of Earth from space, showing the curvature of the planet and city lights at night. The image is a composite of a solid blue background on the left and a satellite photograph of Earth on the right. The Earth's surface is dark, with numerous bright yellow and orange lights representing cities and urban areas. The horizon of the Earth is visible as a curved line separating the dark surface from the deep blue of space.

Section 3

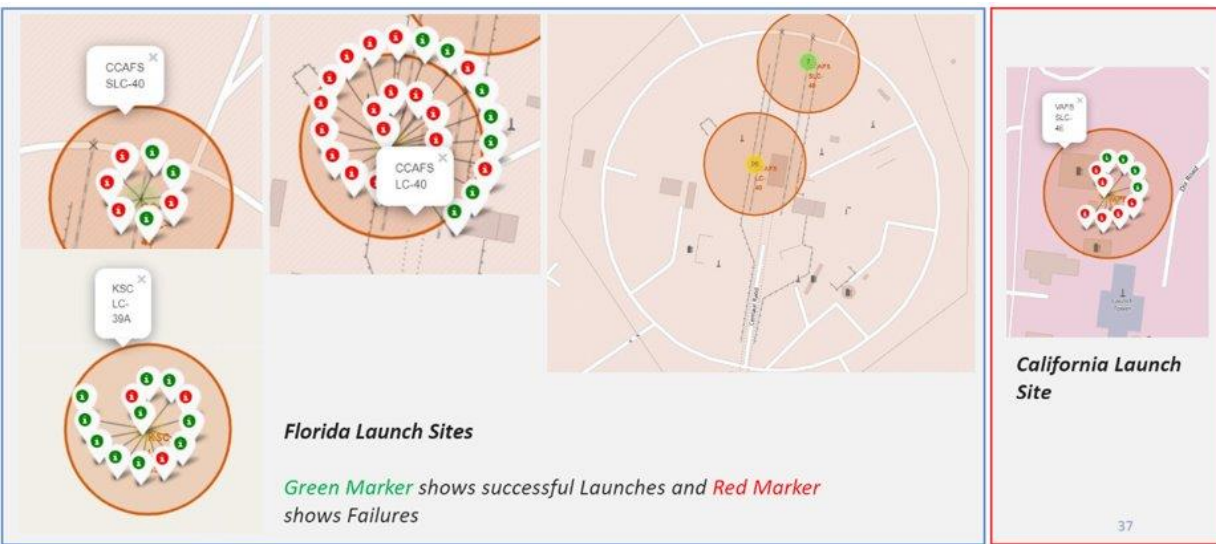
# Launch Sites Proximities Analysis

# All launch sites global map markers

---

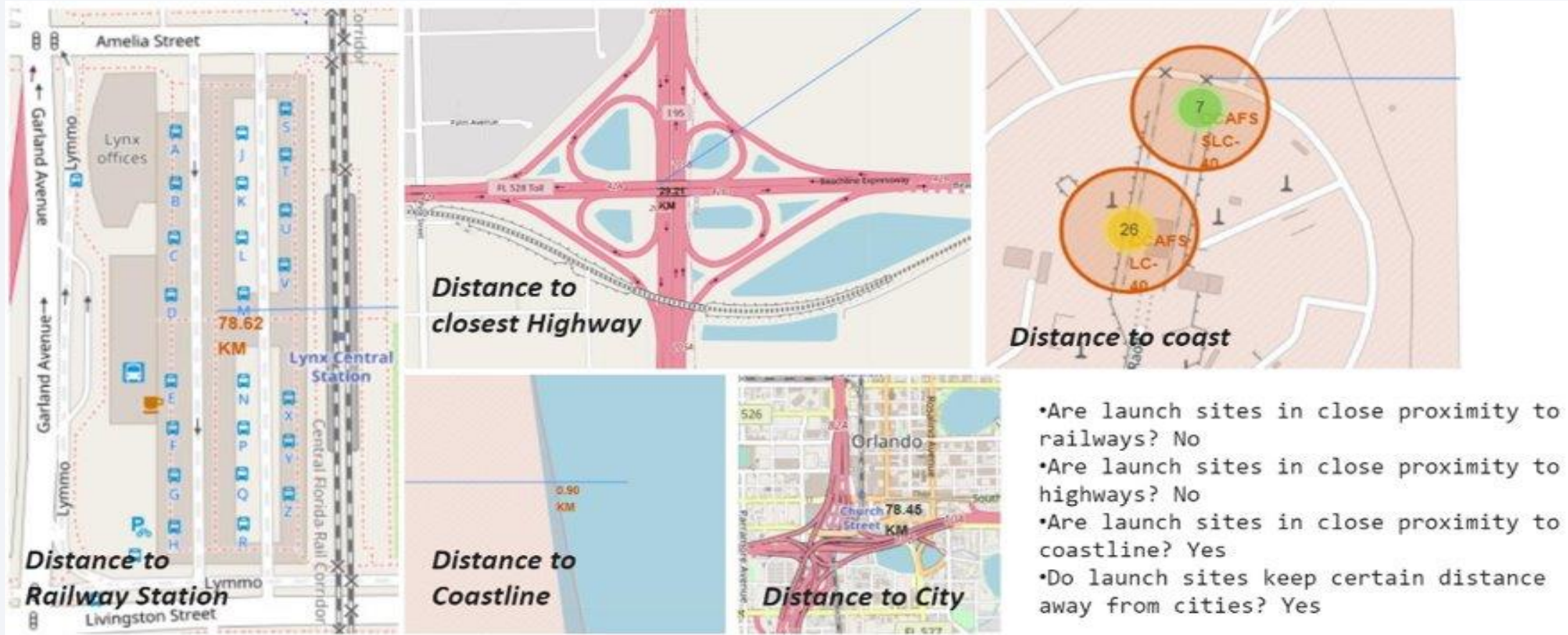


# Markers showing launch sites with color labels





# Launch Site distance to landmarks





Section 4

# Build a Dashboard with Plotly Dash



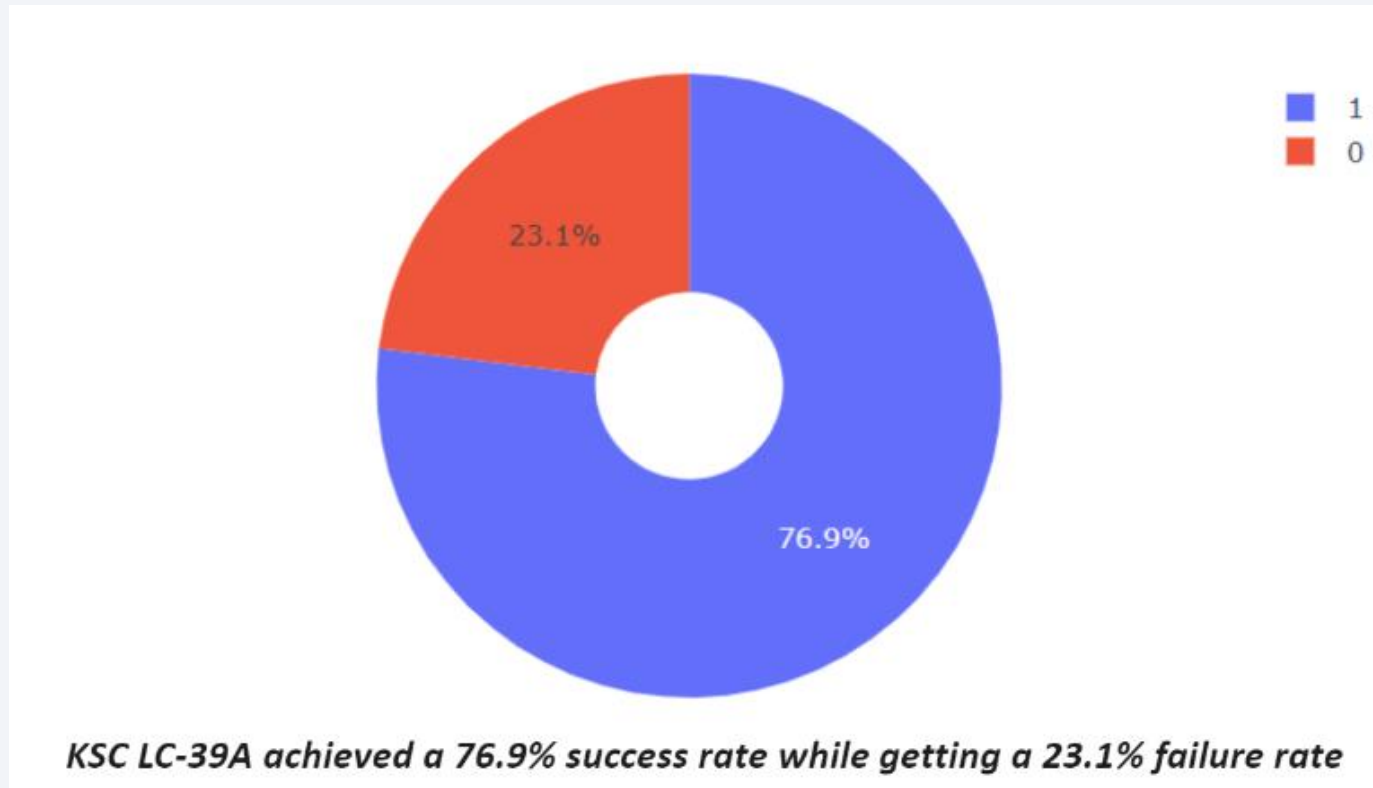
## Pie chart showing the success percentage achieved by each launch site

---



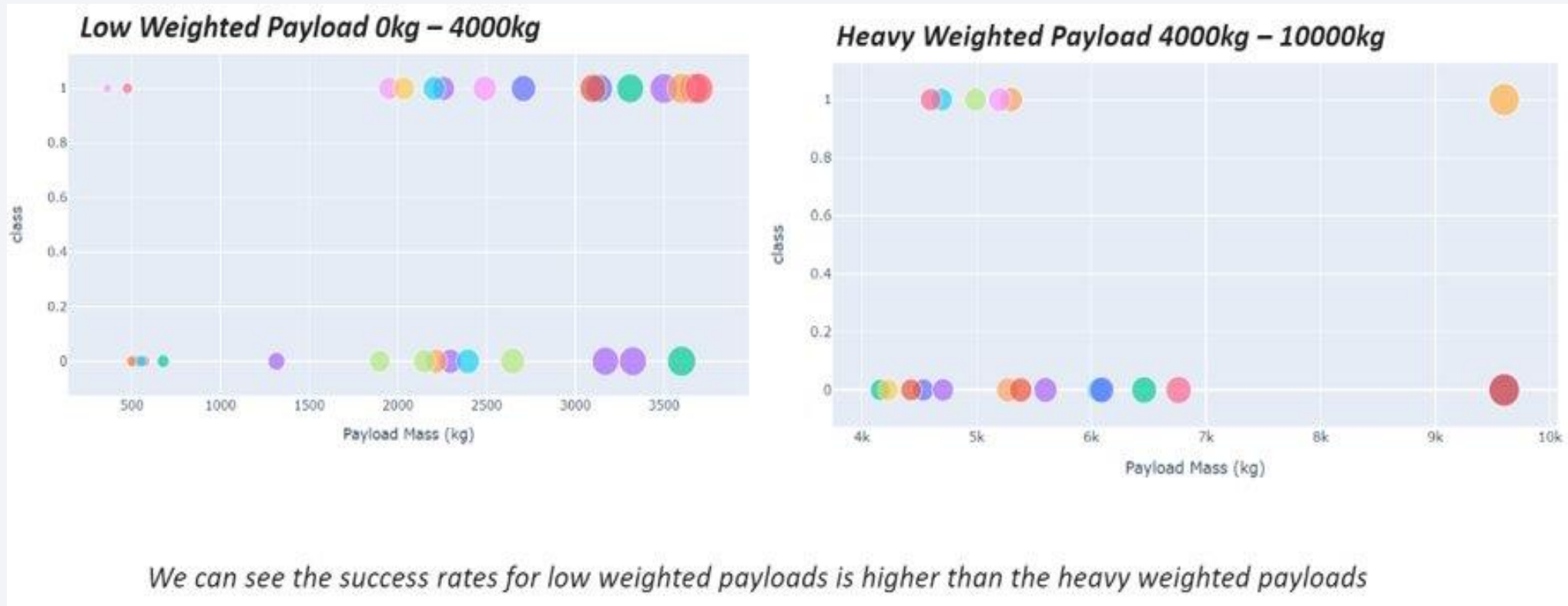
## Pie chart showing the Launch site with the highest launch success ratio

---



## Scatter plot of Payload vs Launch Outcome for all sites, with different payload selected in the range slider

---



Section 5

# Predictive Analysis (Classification)



# Classification Accuracy

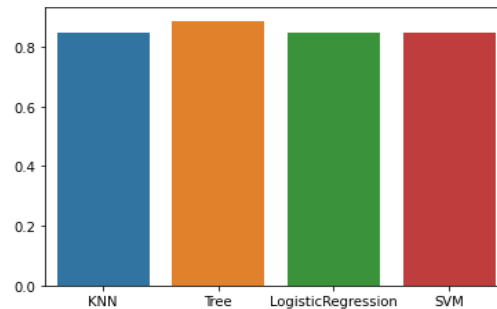
The decision tree classifier is the model with the highest classification accuracy

```
In [38]: algorithms = {'KNN':knn_cv.best_score_, 'Tree':tree_cv.best_score_, 'LogisticRegression':logreg_cv.best_score_, 'SVM':svm_cv.best_score_}
bestalgorithm = max(algorithms, key=algorithms.get)
print('Best Algorithm is',bestalgorithm,'with a score of',algorithms[bestalgorithm])
if bestalgorithm == 'Tree':
    print('Best Params is :',tree_cv.best_params_)
if bestalgorithm == 'KNN':
    print('Best Params is :',knn_cv.best_params_)
if bestalgorithm == 'LogisticRegression':
    print('Best Params is :',logreg_cv.best_params_)

import seaborn as sns
sns.barplot(list(algorithms.keys()), list(algorithms.values()))

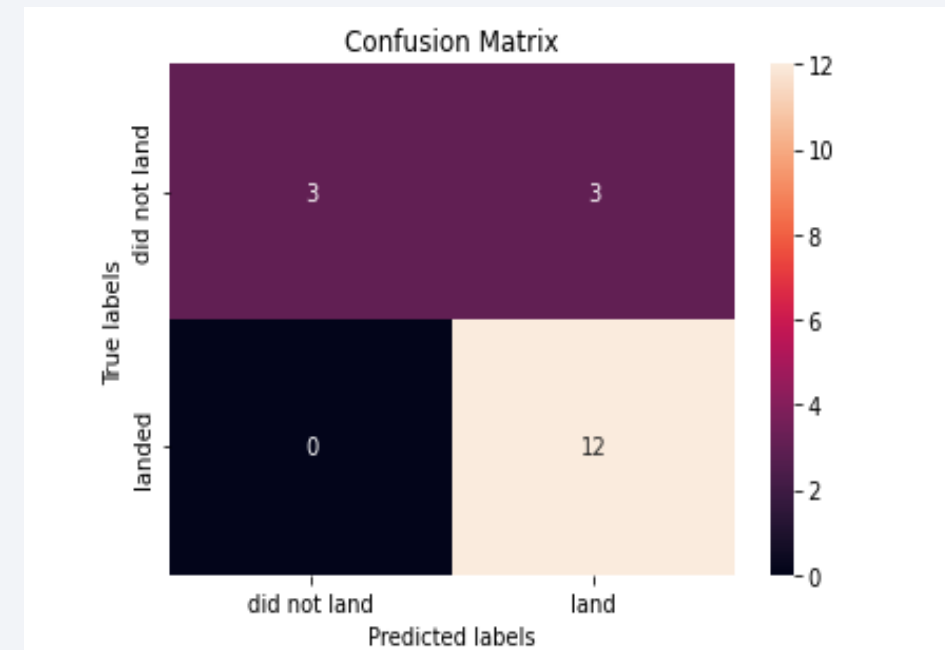
Best Algorithm is Tree with a score of 0.8875
Best Params is : {'criterion': 'gini', 'max_depth': 10, 'max_features': 'auto', 'min_samples_leaf': 4, 'min_samples_split': 2, 'splitter': 'best'}
/opt/conda/envs/Python-3.9/lib/python3.9/site-packages/seaborn/_decorators.py:36: FutureWarning: Pass the following variables as keyword args: x, y. From
m version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or mi
sinterpretation.
  warnings.warn(
<AxesSubplot:>
```

Out[38]:



# 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!

