



IBM Developer  
SKILLS NETWORK

# Winning Space Race with Data Science

Chen  
02/2023



# Outline

---

- Executive Summary
- Introduction
- Methodology
- Results
- Conclusion
- Appendix

# 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 visualization
  - Interactive visual analytics with Folium
  - Machine learning prediction
- Summary of all results
  - Exploratory data analysis results
  - Interactive analytics in screenshots
  - Predictive analysis results

# 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 alternative 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, and evaluate classification models

# 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 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 the data, clean the requested data and did some basic data wrangling and formatting
- The link of the notebook is <https://github.com/nehcgnak/CS-DS-Learning/blob/b3c73b5e5f2ac550427a9504c96399b412c34f20/jupyter-labs-spacex-data-collection-api.ipynb>

Task 1: Request and parse the SpaceX launch data using the GET request

To make the requested JSON results more consistent, we will use the following static response object for this project:

```
static_json_url='https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBM-DS0321EN-SkillsNetwork/datasets/API_call_spacex_api.json'
```

We should see that the request was successful with the 200 status response code

```
response.status_code
```

200

Now we decode the response content as a Json using `.json()` and turn it into a Pandas dataframe using `.json_normalize()`

```
# Use json_normalize meethod to convert the json result into a dataframe
data = pd.json_normalize(response.json())
```

Using the dataframe `data` print the first 5 rows

```
# Get the head of the dataframe
data.head()
```



# Data Collection - Scraping

- We applied web scrapping to collect 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/nehcgnaK/CS-DS-Learning/blob/b3c73b5e5f2ac550427a9504c96399b412c34f20/jupyter-labs-webscraping.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.

```
# use requests.get() method with the provided static_url
# assign the response to a object
response = requests.get(static_url).text
```

Create a BeautifulSoup object from the HTML response

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

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

```
# Use soup.title attribute
print(soup.title)
```

```
<title>List of Falcon 9 and Falcon Heavy launches - Wikipedia</title>
```

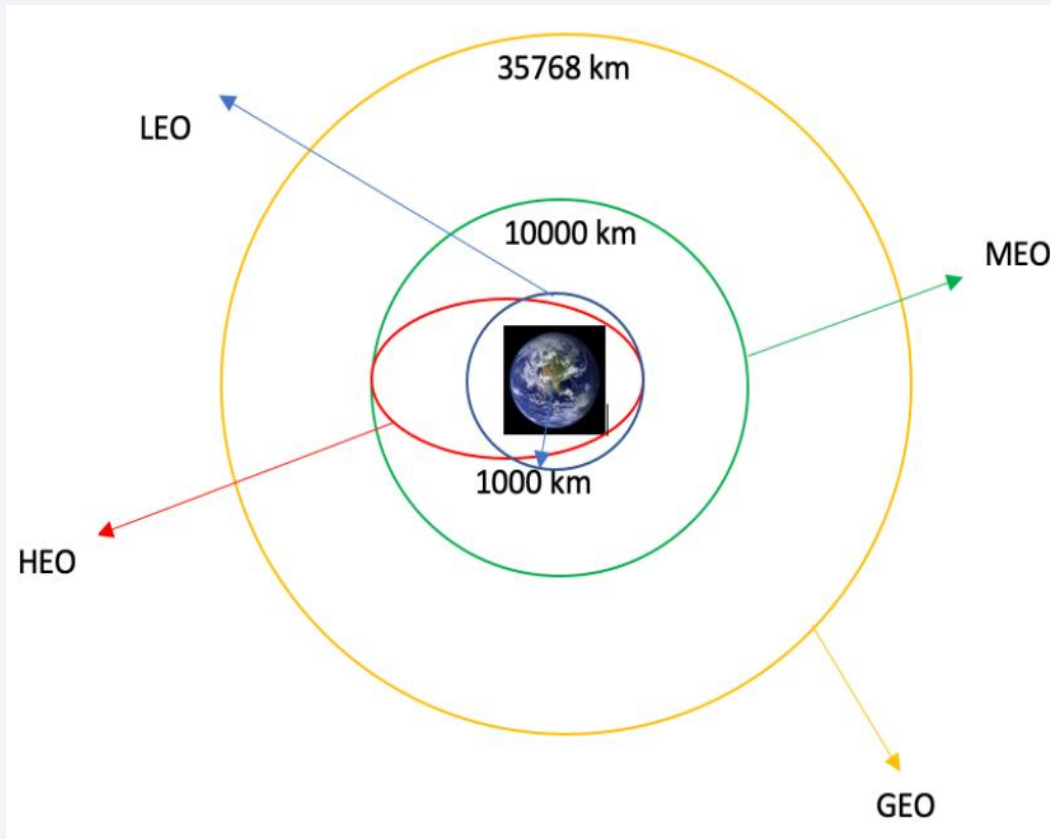
TASK 2: Extract all column/variable names from the HTML table header

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, please check the external reference link towards the end of this lab

```
# Use the find_all function in the BeautifulSoup object, with element type `table`
# Assign the result to a list called `html_tables`
html_tables = soup.find_all("table")
print(html_tables)
```

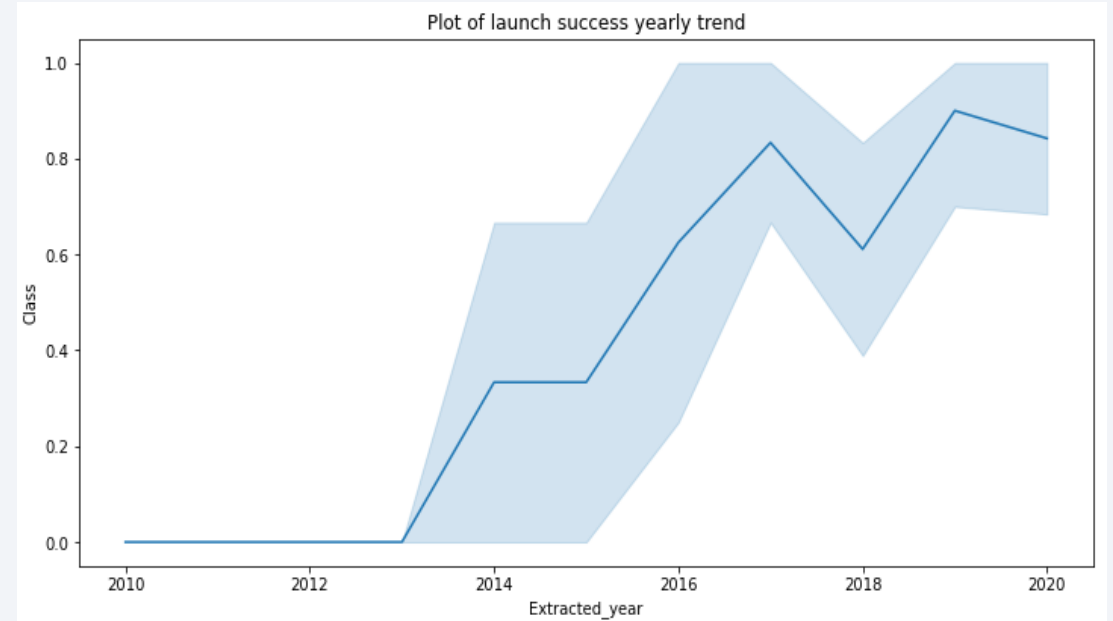
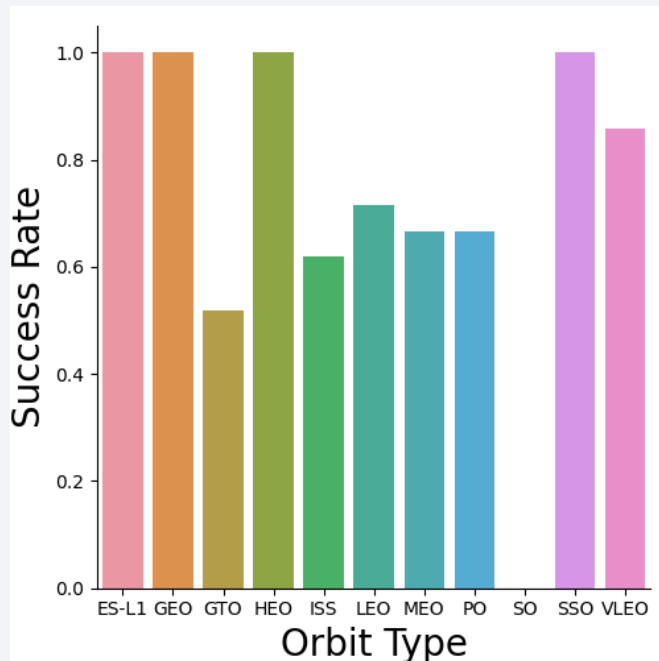
# Data Wrangling



- 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 a csv file
- The link to the notebook is [https://github.com/nehcgnak/CS-DS-Learning/blob/b3c73b5e5f2ac550427a9504c96399b412c34f20/IBM-DS0321EN-SkillsNetwork\\_labs\\_module\\_1\\_L3\\_labs-jupyter-spacex-data\\_wrangling\\_jupyterlite.jupyterlite.ipynb](https://github.com/nehcgnak/CS-DS-Learning/blob/b3c73b5e5f2ac550427a9504c96399b412c34f20/IBM-DS0321EN-SkillsNetwork_labs_module_1_L3_labs-jupyter-spacex-data_wrangling_jupyterlite.jupyterlite.ipynb)

# EDA with Data Visualization

- We explored the data by visualizing the relationship between flight number vs. launch site, payload vs. launch site, success rate of each orbit type, flight number vs orbit type, and yearly trend of the launch success



- The link to the notebook is [https://github.com/nehcgnak/CS-DS-Learning/blob/b3c73b5e5f2ac550427a9504c96399b412c34f20/IBM-DS0321EN-SkillsNetwork\\_labs\\_module\\_2\\_jupyter-labs-eda-dataviz.ipynb.jupyterlite.ipynb](https://github.com/nehcgnak/CS-DS-Learning/blob/b3c73b5e5f2ac550427a9504c96399b412c34f20/IBM-DS0321EN-SkillsNetwork_labs_module_2_jupyter-labs-eda-dataviz.ipynb.jupyterlite.ipynb)

# EDA with SQL

---

- 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/nehcgnak/CS-DS-Learning/blob/b3c73b5e5f2ac550427a9504c96399b412c34f20/jupyter-labs-eda-sql-coursera\\_sqlite.ipynb](https://github.com/nehcgnak/CS-DS-Learning/blob/b3c73b5e5f2ac550427a9504c96399b412c34f20/jupyter-labs-eda-sql-coursera_sqlite.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 different booster versions

# Predictive Analysis (Classification)

---

- We loaded the data using numpy and pandas, transformed the data, split the data into training and testing sets.
- We built different machine learning models and tune different hyper-parameters using GridSearchCV
- We used accuracy as the metric for our model, improved the model using feature engineering and hyper-parameter tuning
- We found the best performing classification model
- The link to the notebook is [https://github.com/nehcgnak/CS-DS-Learning/blob/b3c73b5e5f2ac550427a9504c96399b412c34f20/IBM-DS0321EN-SkillsNetwork\\_labs\\_module\\_4\\_SpaceX\\_Machine\\_Learning\\_Prediction\\_Part\\_5.jupyterlite.ipynb](https://github.com/nehcgnak/CS-DS-Learning/blob/b3c73b5e5f2ac550427a9504c96399b412c34f20/IBM-DS0321EN-SkillsNetwork_labs_module_4_SpaceX_Machine_Learning_Prediction_Part_5.jupyterlite.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, red, and cyan on the right. These streaks are layered over a faint, grid-like pattern, creating a sense of depth and movement, reminiscent of a digital or data visualization theme.

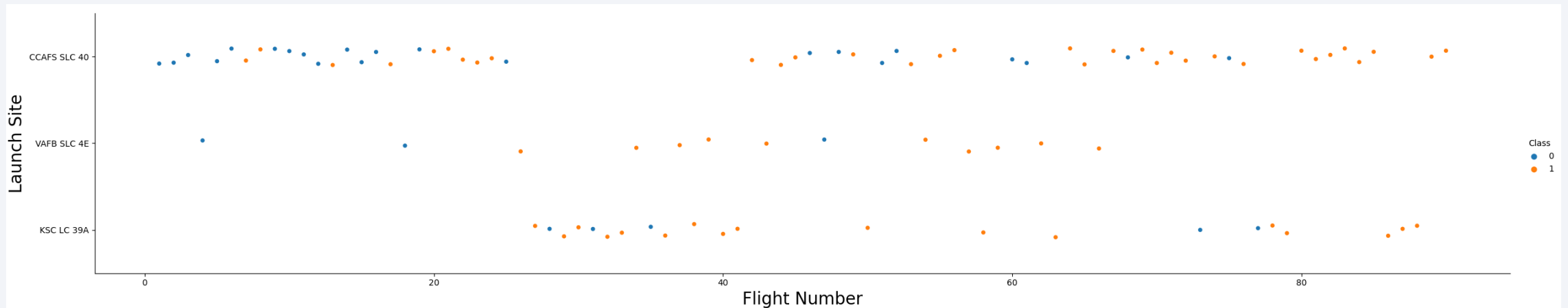
Section 2

# 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 that launch site.

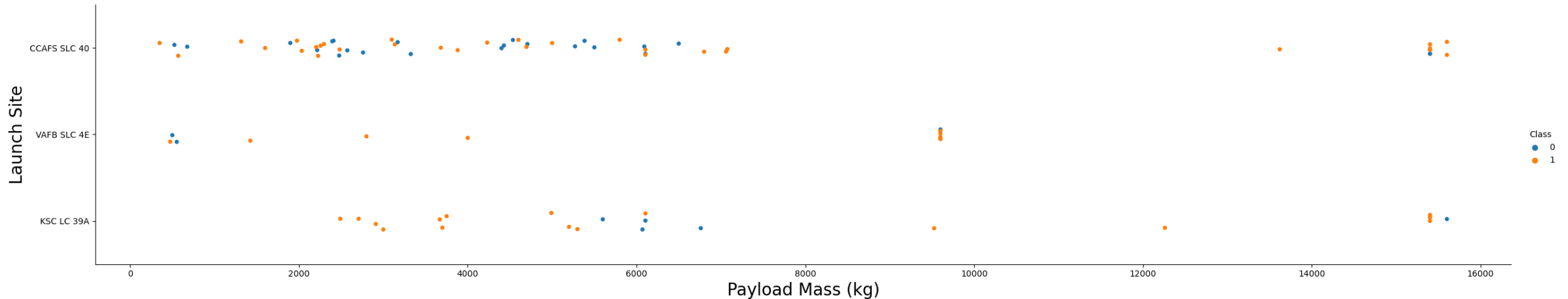




# Payload vs. Launch Site

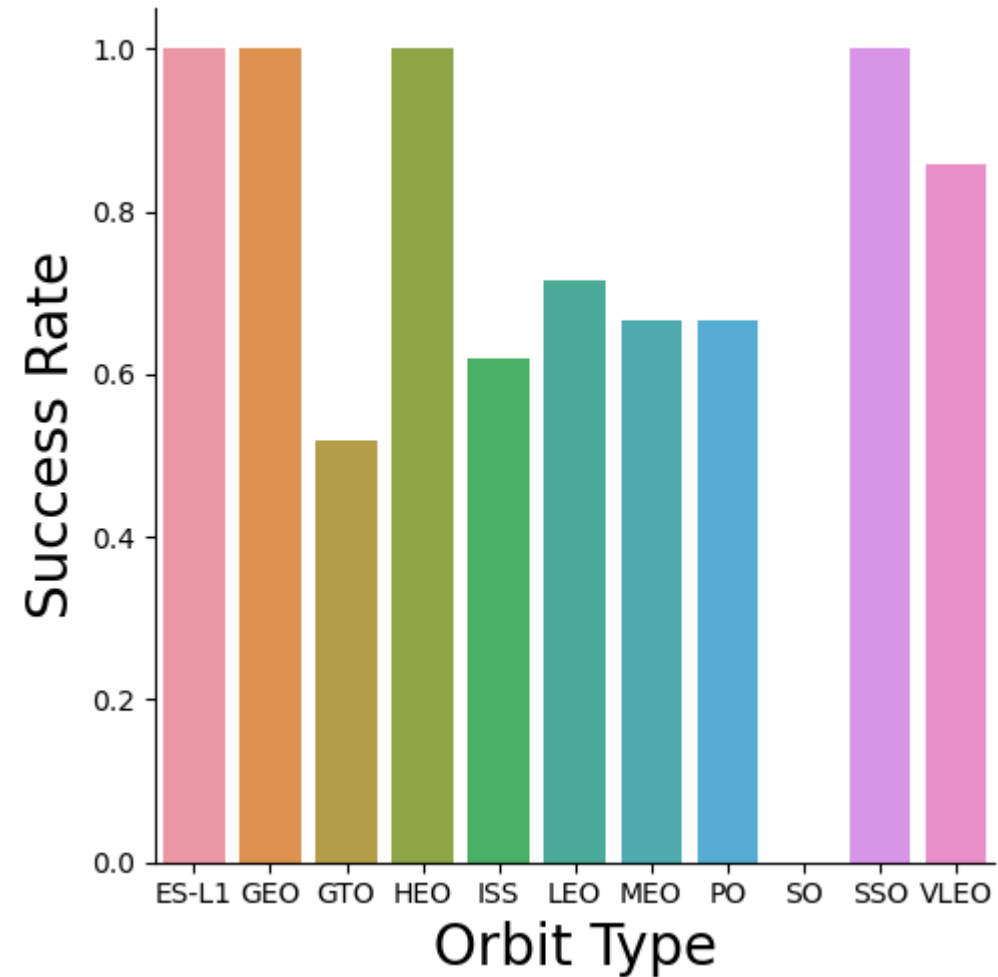


The greater the payload mass for launch site CCAFS SLC 40 the higher the success rate for the rocket.



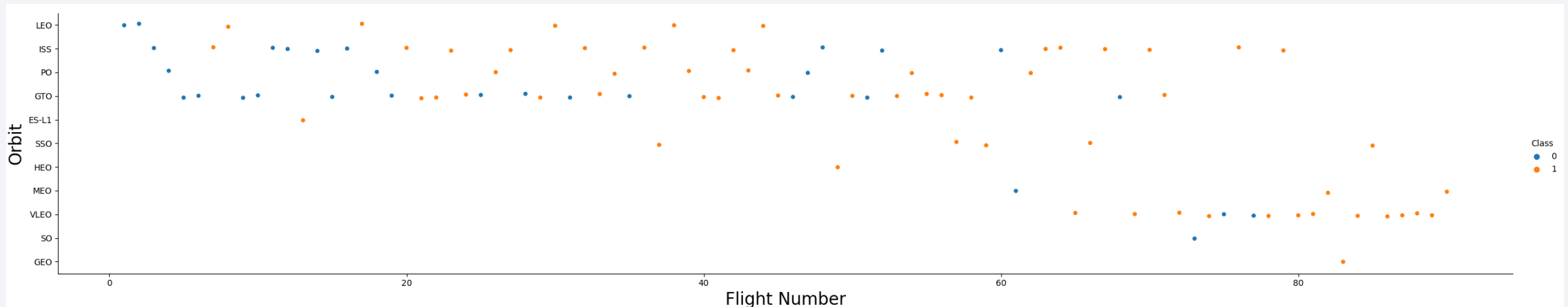
# Success Rate vs. Orbit Type

- From the plot, we can see that ES-L1, GEO, HEO, SSO, VLEO have the higher success rates



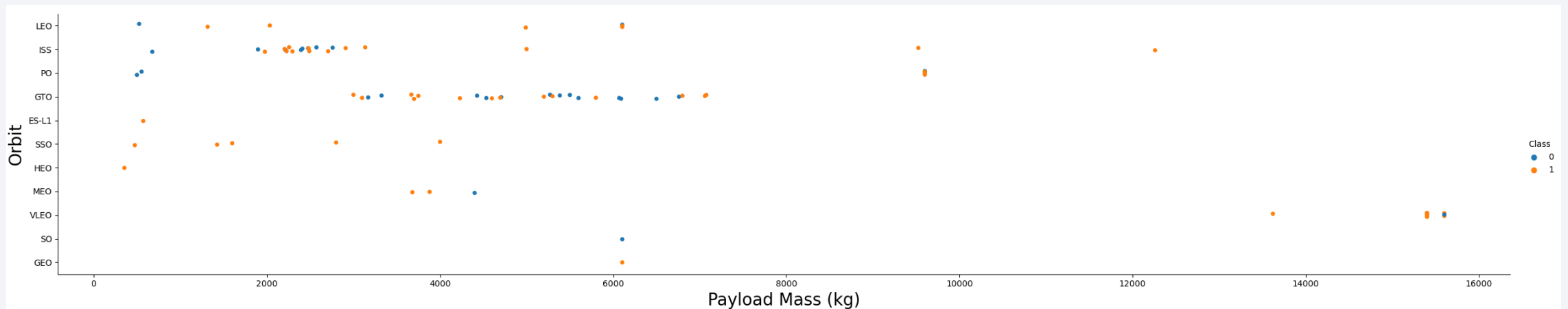
# Flight Number vs. Orbit Type

- The plot below shows the Flight Number vs. Orbit type. We observed 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.



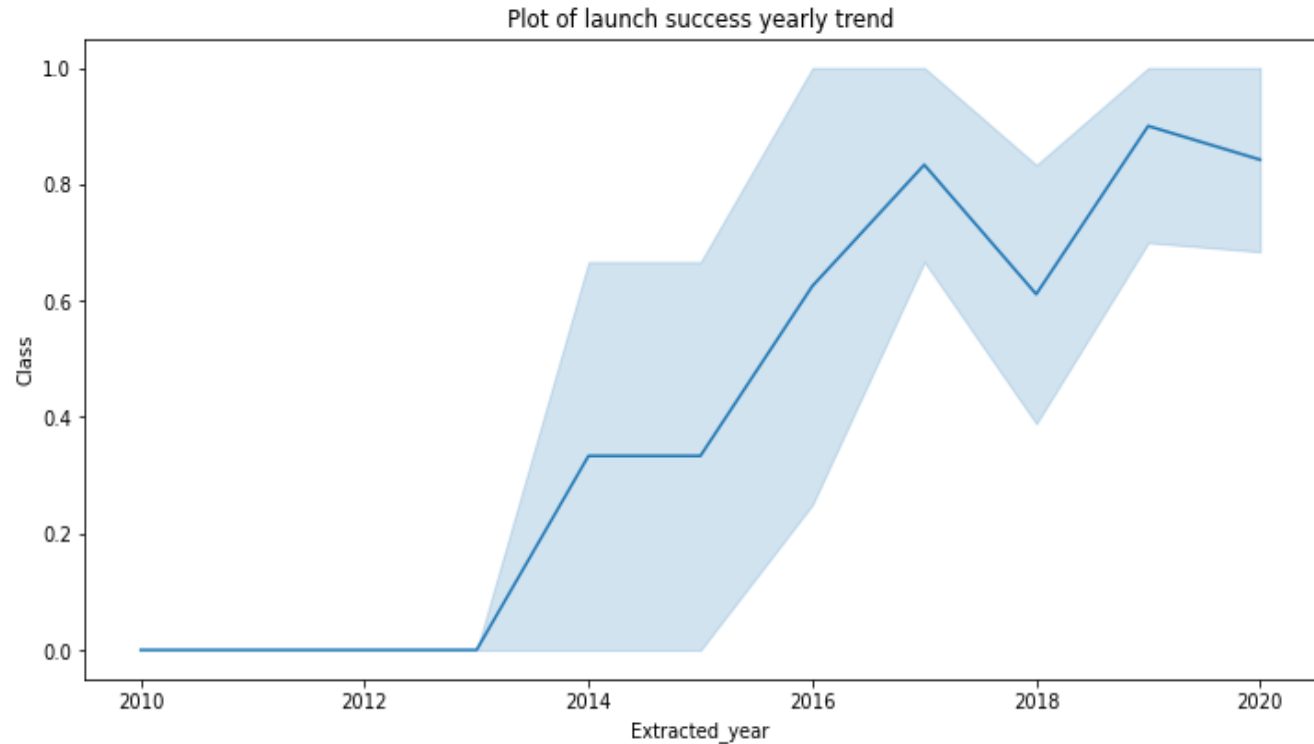
# Payload vs. Orbit Type

- We observed that with heavy payloads, the successful landing were more for PO, LEO, and ISS orbits.



# Launch Success Yearly Trend

- From the plot, we observed that success rate kept on increasing from 2013 to 2020





# All Launch Site Names

- We used the key word **DISTINCT** to show 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 SPACEXTBL
```

```
* 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'

## Task 2

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

```
%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
04-06-2010	18:45:00	F9 v1.0 B0003	CCAFS LC-40	Dragon Spacecraft Qualification Unit	0	LEO	SpaceX	Success	Failure (parachute)
08-12-2010	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)
22-05-2012	07:44:00	F9 v1.0 B0005	CCAFS LC-40	Dragon demo flight C2	525	LEO (ISS)	NASA (COTS)	Success	No attempt
08-10-2012	00:35:00	F9 v1.0 B0006	CCAFS LC-40	SpaceX CRS-1	500	LEO (ISS)	NASA (CRS)	Success	No attempt
01-03-2013	15:10:00	F9 v1.0 B0007	CCAFS LC-40	SpaceX CRS-2	677	LEO (ISS)	NASA (CRS)	Success	No attempt

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

# Total Payload Mass

---

- We calculated the total payload carried by boosters from NASA is 45,596 using the query below

## Task 3

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

```
%sql select sum(PAYLOAD_MASS__KG_) from SPACEXTBL where CUSTOMER = 'NASA (CRS)'
```

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

```
Done.
```

```
sum(PAYLOAD_MASS__KG_)
```

```
45596
```

# Average Payload Mass by F9 v1.1

- We calculated the average payload mass carried by booster version F9 v1.1 is 2,928.4

## Task 4

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

```
%sql select avg(PAYLOAD_MASS_KG_) from SPACEXTBL where BOOSTER_VERSION = 'F9 v1.1'
```

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

```
Done.
```

```
avg(PAYLOAD_MASS_KG_)
```

```
2928.4
```

# First Successful Ground Landing Date

- We observed that the date of the first successful landing outcome on ground pad was 12/22/2015

In [14]:

```
task_5 = '''
    SELECT MIN(Date) AS FirstSuccessfull_landing_date
    FROM SpaceX
    WHERE LandingOutcome LIKE 'Success (ground pad)'
    '''

create_pandas_df(task_5, database=conn)
```

Out[14]:

	firstsuccessfull_landing_date
0	2015-12-22



# Successful Drone Ship Landing with Payload between 4,000 and 6,000

```
In [15]: task_6 = '''
          SELECT BoosterVersion
          FROM SpaceX
          WHERE LandingOutcome = 'Success (drone ship)'
             AND PayloadMassKG > 4000
             AND PayloadMassKG < 6000
          ...
          create_pandas_df(task_6, database=conn)
```

```
Out[15]:
```

	boosterversion
0	F9 FT B1022
1	F9 FT B1026
2	F9 FT B1021.2
3	F9 FT B1031.2

- We used the **WHERE** clause to filter for boosters which had successfully landed on drone ship and applied the **AND** condition to determine successful landing with payload mass greater than 4,000 but less than 6,000 kg

# Total Number of Successful and Failure Mission Outcomes

- We used wildcard “like” and “%” to filter for **WHERE** the mission outcome was a success or a failure

## Task 7

List the total number of successful and failure mission outcomes

```
%sql select mission_outcome, count(*) as total_number from SPACEXTBL 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

- We determined the booster that had carried the maximum payload by 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

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

```
* 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
---------------

# 2015 Launch Records

---

- We used a combination 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

```
List the failed landing_outcomes in drone ship, their booster versions, and launch site names for in year 2015

In [18]: task_9 = '''
          SELECT BoosterVersion, LaunchSite, LandingOutcome
          FROM SpaceX
          WHERE LandingOutcome LIKE 'Failure (drone ship)'
             AND Date BETWEEN '2015-01-01' AND '2015-12-31'
          ...
          create_pandas_df(task_9, database=conn)

Out[18]:
```

	boosterversion	launchsite	landingoutcome
0	F9 v1.1 B1012	CCAFS LC-40	Failure (drone ship)
1	F9 v1.1 B1015	CCAFS LC-40	Failure (drone ship)

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

Rank the count of landing outcomes (such as Failure (drone ship) or Success (ground pad))

```
In [19]: task_10 = '''
          SELECT LandingOutcome, COUNT(LandingOutcome)
          FROM SpaceX
          WHERE DATE BETWEEN '2010-06-04' AND '2017-03-20'
          GROUP BY LandingOutcome
          ORDER BY COUNT(LandingOutcome) DESC
          '''
          create_pandas_df(task_10, database=conn)
```

```
Out[19]:
```

	landingoutcome	count
0	No attempt	10
1	Success (drone ship)	6
2	Failure (drone ship)	5
3	Success (ground pad)	5
4	Controlled (ocean)	3
5	Uncontrolled (ocean)	2
6	Precluded (drone ship)	1
7	Failure (parachute)	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** 06/04/2010 and 03/20/2017
- We applied the **GROUP BY** clause to group the landing outcomes and the **ORDER BY** clause to order the grouped landing outcome in a descending order

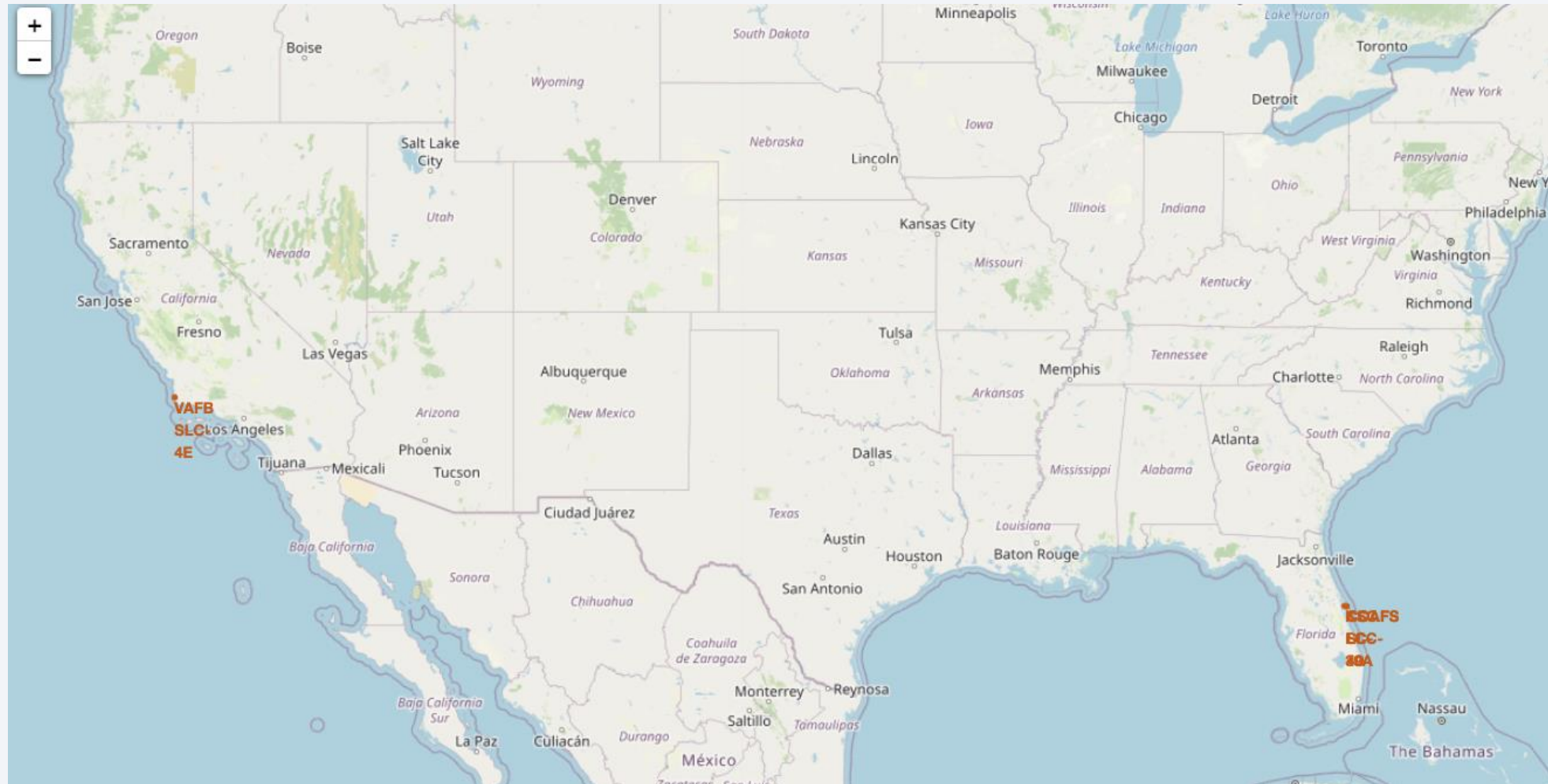
Section 4

# Launch Sites Proximities Analysis





# All Launch Sites Global Map Markers





# Markers Showing Launch Sites with Color Labels



# Launch Site Distance to Landmarks



- Are launch sites in close proximity to railways? No
- Are launch sites in close proximity to highways? No
- Are launch sites in close proximity to coastline? Yes
- Do launch sites keep certain distance away from cities? Yes



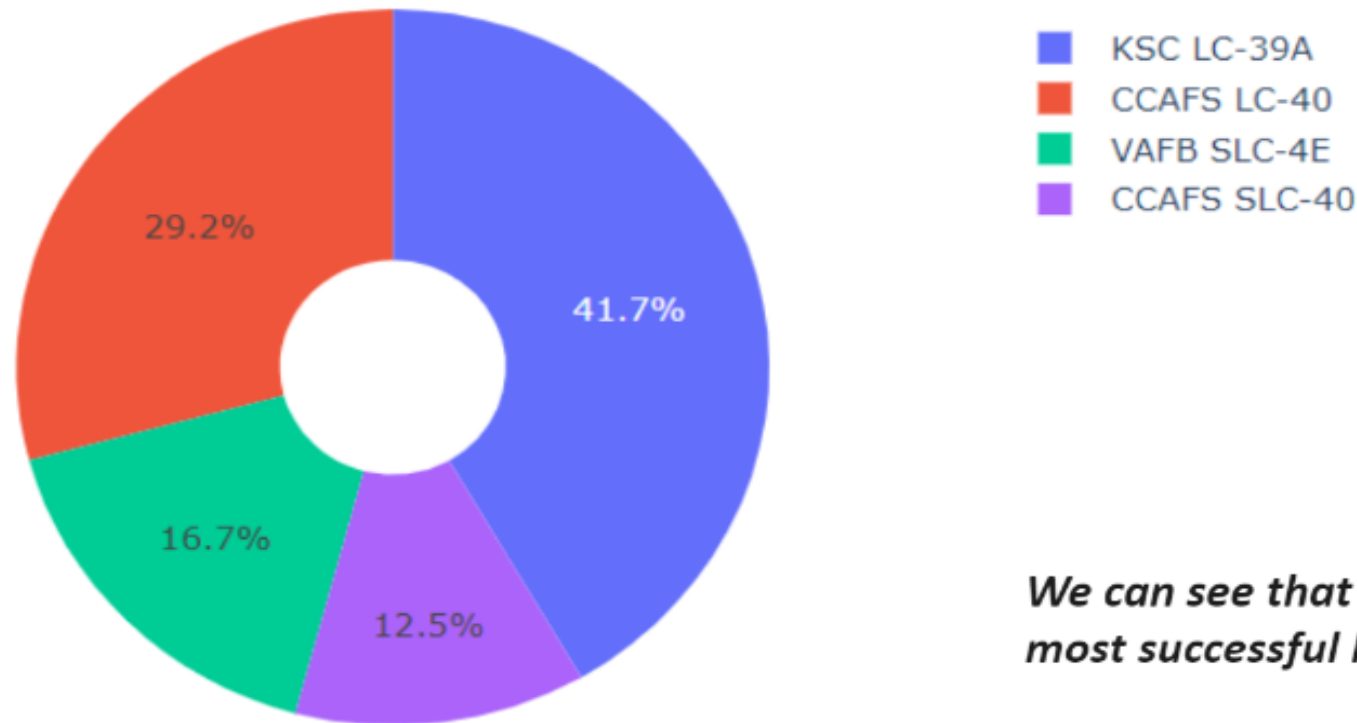


Section 5

# Build a Dashboard with Plotly Dash

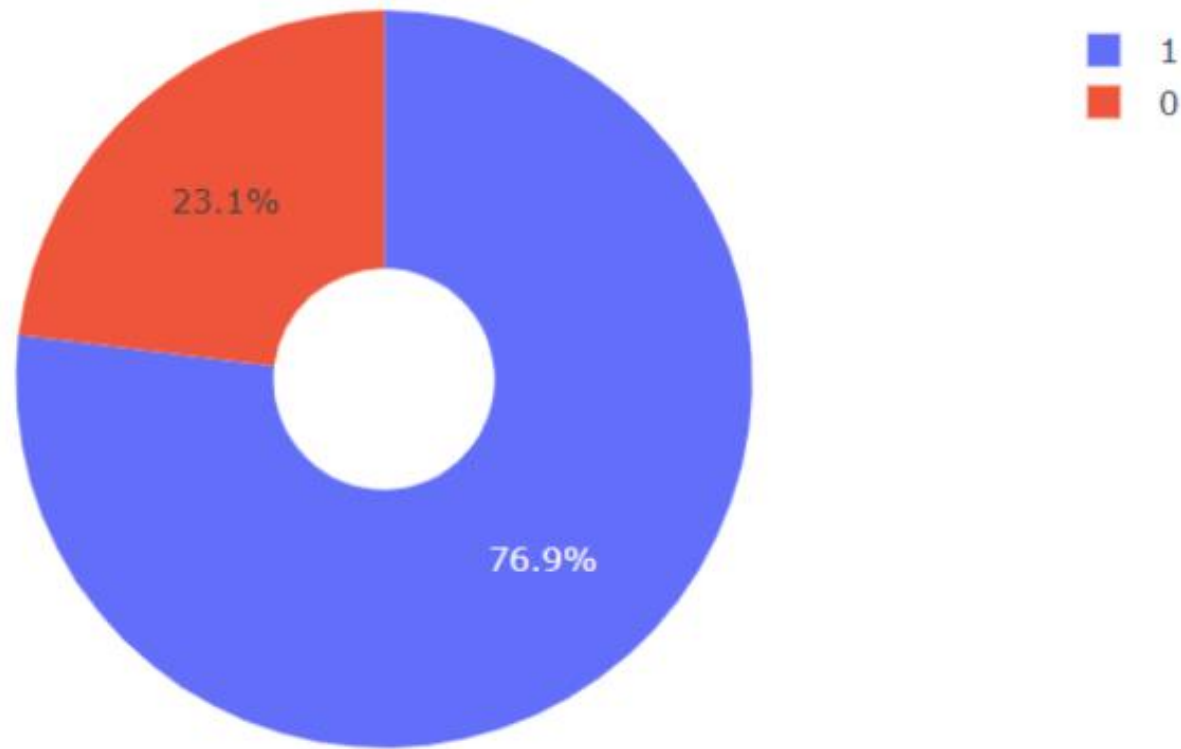
# Pie Chart Showing the Success Percentage Achieved by Each Launch Site

Total Success Launches By all sites



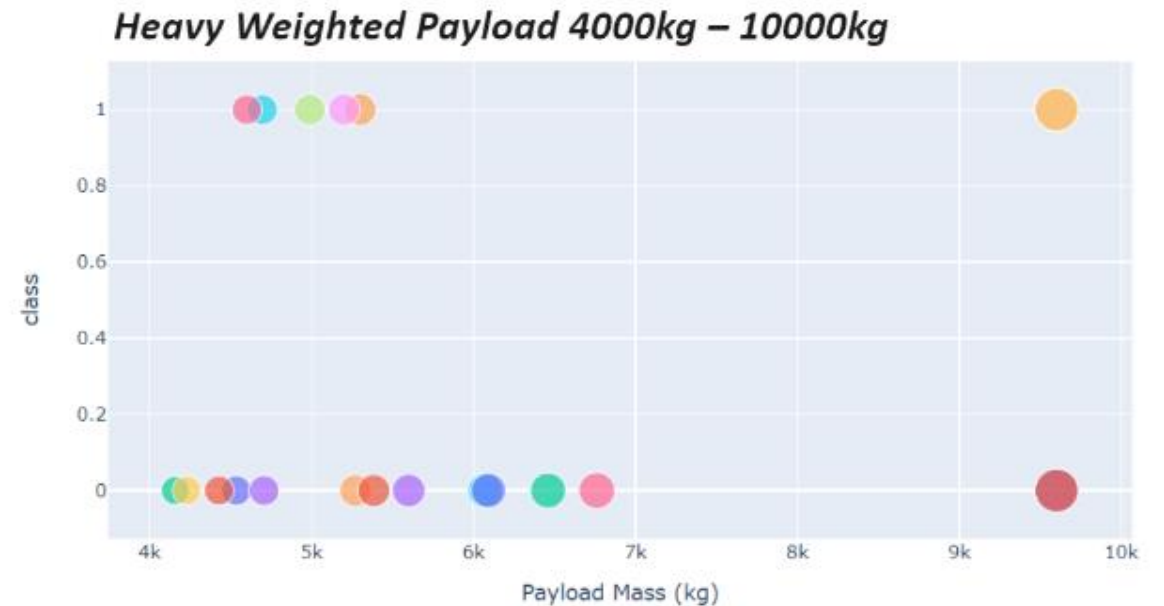
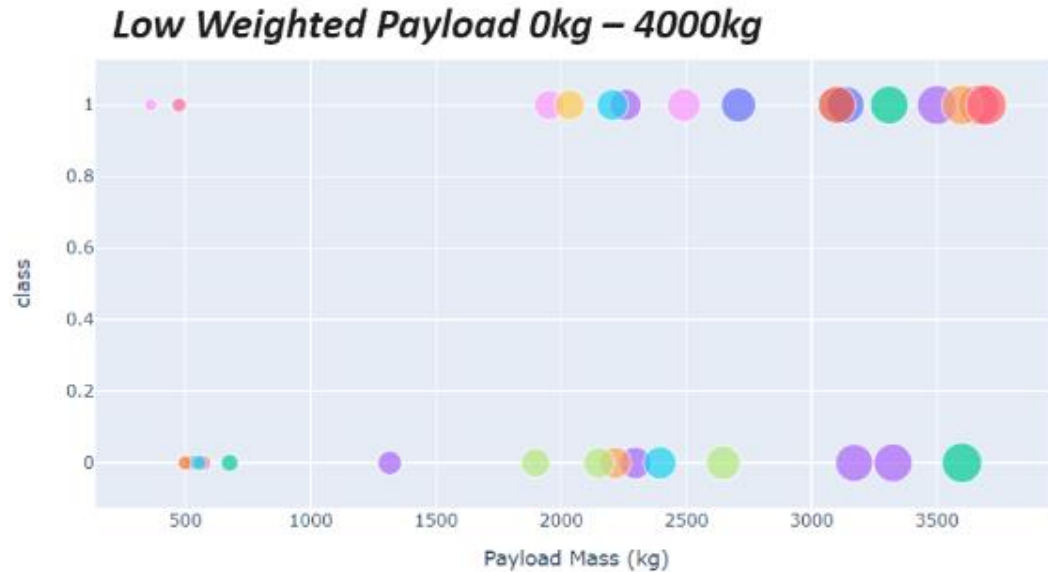
***We can see that KSC LC-39A had the most successful launches from all the sites***

# Pie Chart Showing the Launch Site with the Highest Launch Success Rate



***KSC LC-39A achieved a 76.9% success rate while getting a 23.1% failure rate***

# Scatter Plot of Payload vs. Launch Outcome for All Sites, with Different Payload Selected in the Range Slider



*We can see the success rates for low weighted payloads is higher than the heavy weighted payloads*



Section 6

# Predictive Analysis (Classification)

# Classification Accuracy

- All machine learning classifiers have the same accuracy in the testing dataset

## TASK 12

Find the method performs best:

```
print("Model\t\tAccuracy\tTestAccuracy")#, logreg_cv.best_score_)
print("LogReg\t\t{}\t\t{}".format((logreg_cv.best_score_).round(5), logreg_cv.score(X_test, Y_test).round(5)))
print("SVM\t\t{}\t\t{}".format((svm_cv.best_score_).round(5), svm_cv.score(X_test, Y_test).round(5)))
print("Tree\t\t{}\t\t{}".format((tree_cv.best_score_).round(5), tree_cv.score(X_test, Y_test).round(5)))
print("KNN\t\t{}\t\t{}".format((knn_cv.best_score_).round(5), knn_cv.score(X_test, Y_test).round(5)))

comparison = {}

comparison['LogReg'] = {'Accuracy': logreg_cv.best_score_.round(5), 'TestAccuracy': logreg_cv.score(X_test, Y_test).round(5)}
comparison['SVM'] = {'Accuracy': svm_cv.best_score_.round(5), 'TestAccuracy': svm_cv.score(X_test, Y_test).round(5)}
comparison['Tree'] = {'Accuracy': tree_cv.best_score_.round(5), 'TestAccuracy': tree_cv.score(X_test, Y_test).round(5)}
comparison['KNN'] = {'Accuracy': knn_cv.best_score_.round(5), 'TestAccuracy': knn_cv.score(X_test, Y_test).round(5)}
```

Model	Accuracy	TestAccuracy
LogReg	0.84643	0.83333
SVM	0.84821	0.83333
Tree	0.875	0.83333
KNN	0.84821	0.83333



# 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 - 2020
- Orbits ES-L1, GEO, HEO, SSO, VLEO had the most success rate
- KSC LC-39A had the most successful launches in any sites
- All machine learning classifiers have the same accuracy

Thank you!

