**IBM Developer**
**SKILLS NETWORK**

# Winning Space Race with Data Science

Moustafa Abada
07/25/2023

# Outline

- Executive Summary

- Introduction

- Methodology

- Results

- Conclusion

- Appendix

# Executive Summary

- Summary of methodologies

  - Data collection and data wrangling

  - EDA and interactive visual analytics

  - Predictive analysis

  - EDA with visualization results

  - EDA with SQL results

  - interactive map with Folium results

  - Plotly Dash dashboard

  - Predictive analysis

- Summary of all results

  - Exploratory Data Analysis result

  - Interactive analytics in screenshots

  - Predictive Analytics 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

# 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 models

# Data Collection

- The data was collected using various methods

  - Collect data through a get request to the SpaceX API.

  - Decode the response content into Json format using the .json() function and convert it into a pandas dataframe using .json_normalize().

  - Clean the data, check for missing values, and fill in missing values where necessary.

  - Perform web scraping from Wikipedia for Falcon 9 launch records using BeautifulSoup to extract the launch records as an HTML table, parse it, and convert it into a pandas dataframe for future analysis.

# Data Collection – SpaceX API

- To gather the data, we utilized the get request to the SpaceX API. Following that, we performed basic data wrangling and formatting, including data cleaning, to prepare it for further analysis.

- The link to the notebook is https://github.com/moustafa126/IB M-Data-Science-Capstone-SpaceX/blob/main/Data_Collection _Via_SpaceX_API.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

- To obtain Falcon 9 launch records, we utilized web scraping with BeautifulSoup. We then parsed the resulting table and transformed it into a pandas dataframe.

- The link to the notebook is https://github.com/moustafa 126/IBM-Data-Science-Capstone-SpaceX/blob/main/Data_Coll ection_with_Web_Scraping.ip ynb



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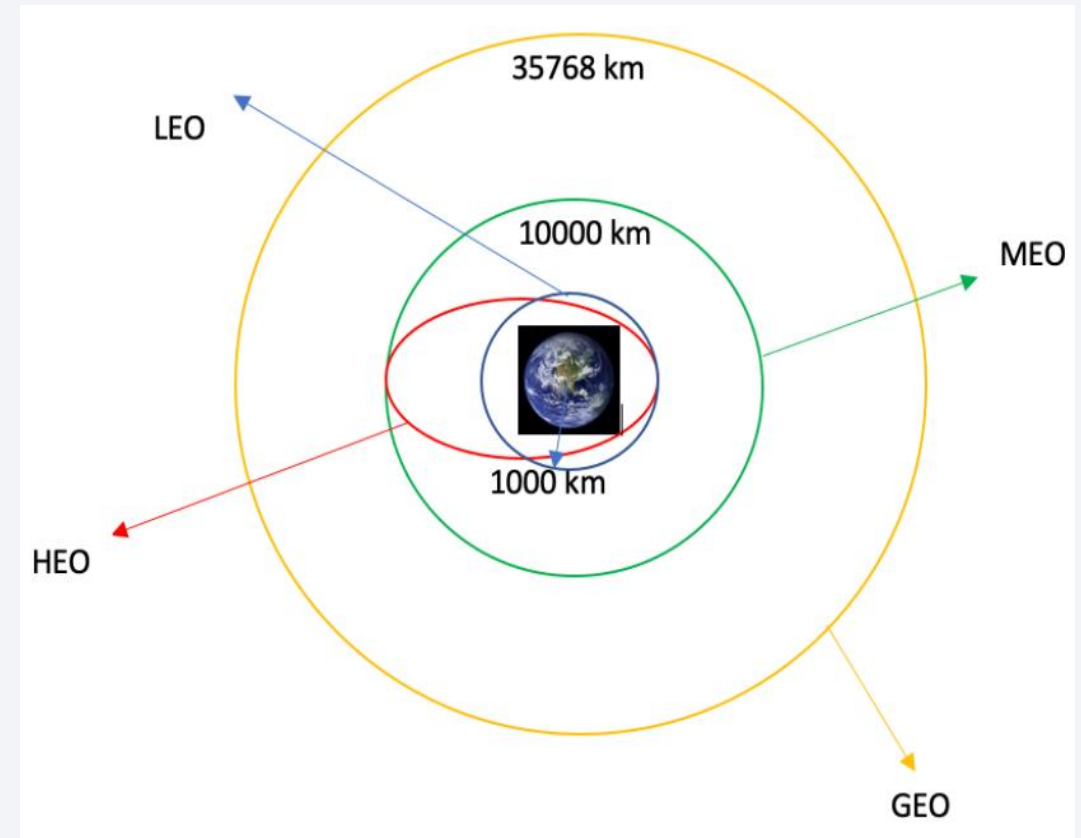
4. Create a dataframe by parsing the launch HTML tables
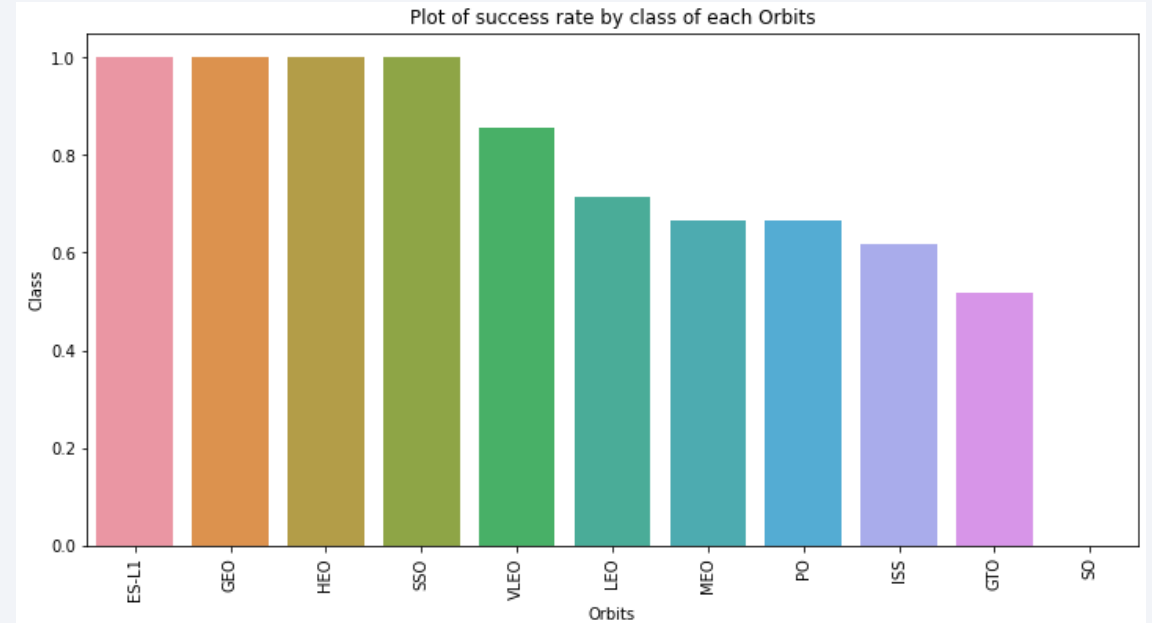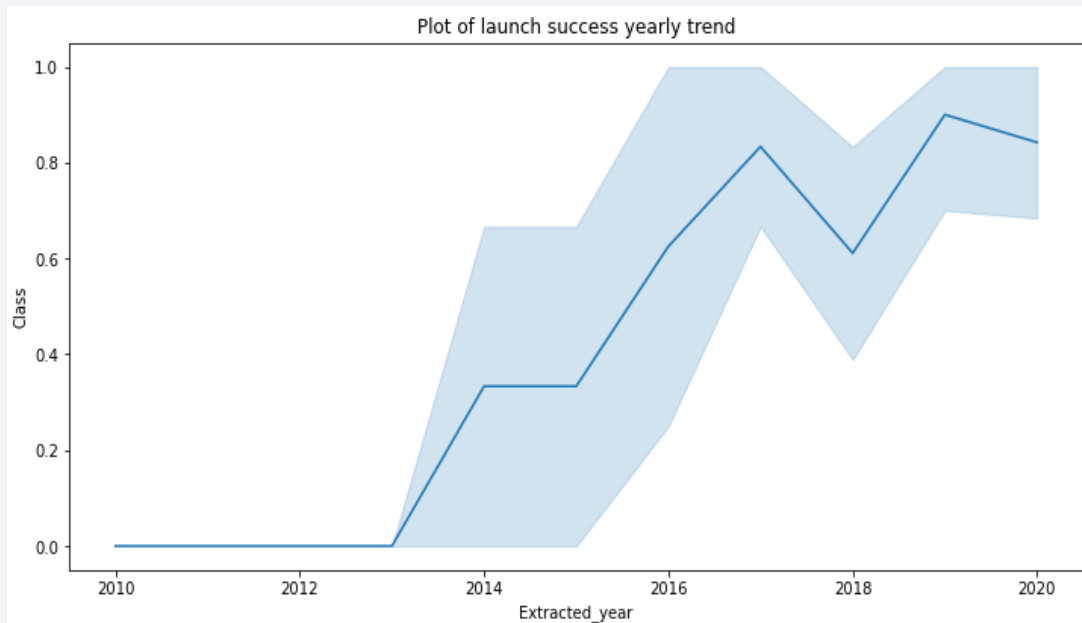5. Export data to csv

# Data Wrangling

- Perform exploratory data analysis to determine the training labels.

- Calculate the number of launches at each site and the frequency and occurrence of each orbit.

- Generate a landing outcome label using data from the outcome column.

- Export the results to a csv file.

- The link to the notebook is https://github.com/moustafa126/IBM-Data-Science-Capstone-SpaceX/blob/main/Data_Wrangling.ipynb

# EDA with Data Visualization

- We analyzed the data by creating visual representations of various relationships. These included the relationship between flight number and launch site, payload and launch site, success rate of each orbit type, flight number and orbit type, and the yearly trend in launch success.

- The link to the notebook is https://github.com/moustafa126/IBM-Data-Science-Capstone-SpaceX/blob/main/EDA_with_Visualization.ipynb



Plot of launch success yearly trend



Plot of success rate by class of each Orbits

# EDA with SQL

- Load the SpaceX dataset into a PostgreSQL database directly from the Jupyter Notebook.

- Use SQL to conduct exploratory data analysis and gain insight from the data.

- Write queries to answer various questions, such as identifying the unique launch sites in the space mission, determining the total payload mass carried by boosters launched by NASA (CRS), calculating the average payload mass carried by booster version F9 v1.1, and determining the total number of successful and failed mission outcomes.

- Use SQL to identify the booster version and launch site names for failed landing outcomes in drone ships.

- The link to the notebook is: https://github.com/moustafa126/IBM-Data-Science-Capstone-SpaceX/blob/main/EDA_with_SQL.ipynb

# Build an Interactive Map with Folium

- Add various map objects, such as markers, circles, and lines, to a folium map to indicate the success or failure of launches at each launch site.

- Assign launch outcomes to class 0 (failure) or 1 (success).

- Use color-labeled marker clusters to identify launch sites with relatively high success rates.

- Calculate distances between launch sites and their proximities to answer questions such as whether launch sites were located near railways, highways, and coastlines, or whether they maintained a certain distance from nearby cities.

- The link to the notebook is: https://github.com/moustafa126/IBM-Data-Science-Capstone-SpaceX/blob/main/Interactive_Visual_Analytics_with_Folium.ipynb

# Build a Dashboard with Plotly Dash

- Build an interactive dashboard using Plotly Dash.

- Plot pie charts displaying the total number of launches at each site.

- Plot scatter plots showing the relationship between the outcome and payload mass (in kg) for different booster versions.

- The link to the notebook is https://github.com/moustafa126/IBM-Data-Science-Capstone-SpaceX/blob/main/Plotly%20Dash.py

# Predictive Analysis (Classification)

- Load the data using NumPy and Pandas libraries, transform it, and split it into a training set and a testing set.

- Build different machine learning models and tune various hyperparameters using GridSearchCV.

- Evaluate the models using accuracy as the performance metric.

- Improve the models using feature engineering and algorithm tuning.

- Identify the best performing classification model among the various models tested.

- The link to the notebook is https://github.com/moustafa126/IBM-Data-Science-Capstone-SpaceX/blob/main/Machine_Learning_Prediction.ipynb

# Results

- Exploratory data analysis results

- Interactive analytics demo in screenshots

- Predictive analysis results

Section 2
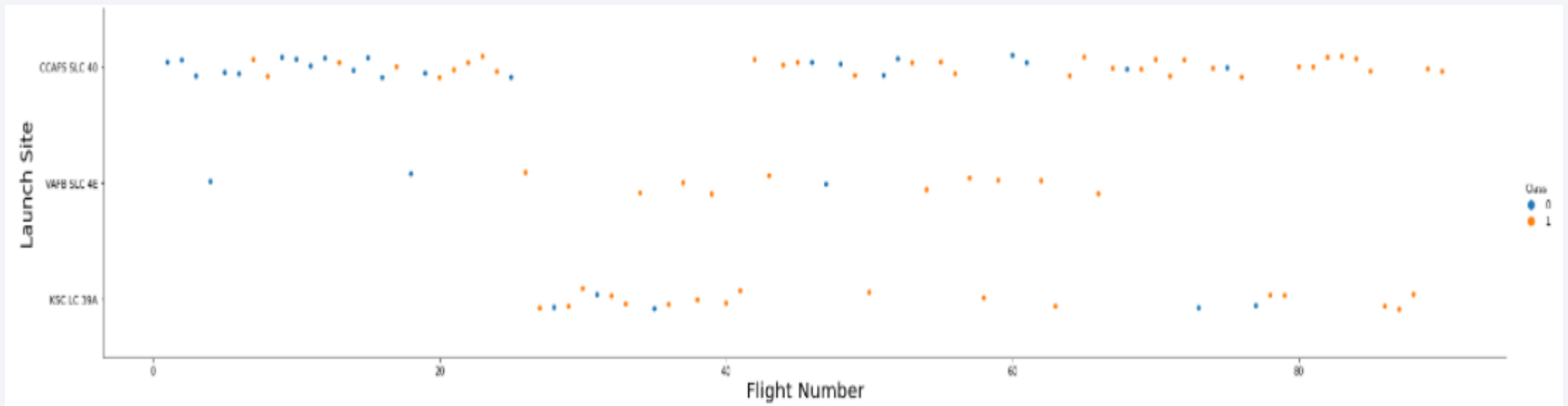
# Insights drawn from EDA

# 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.
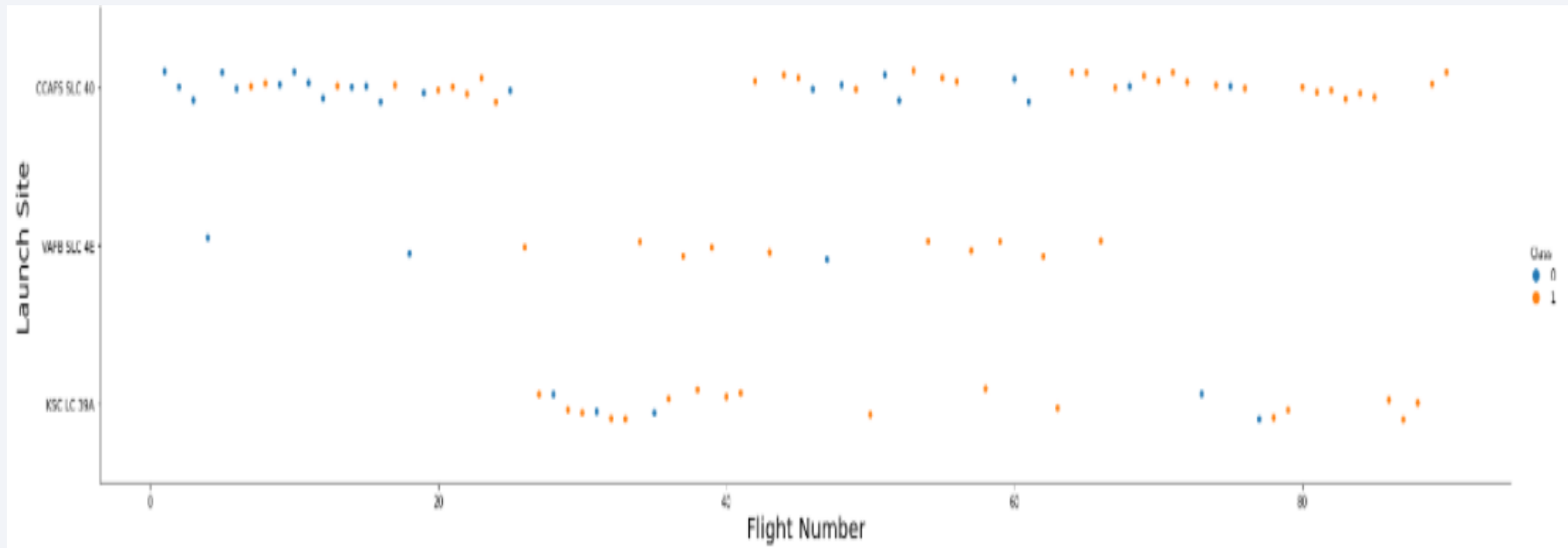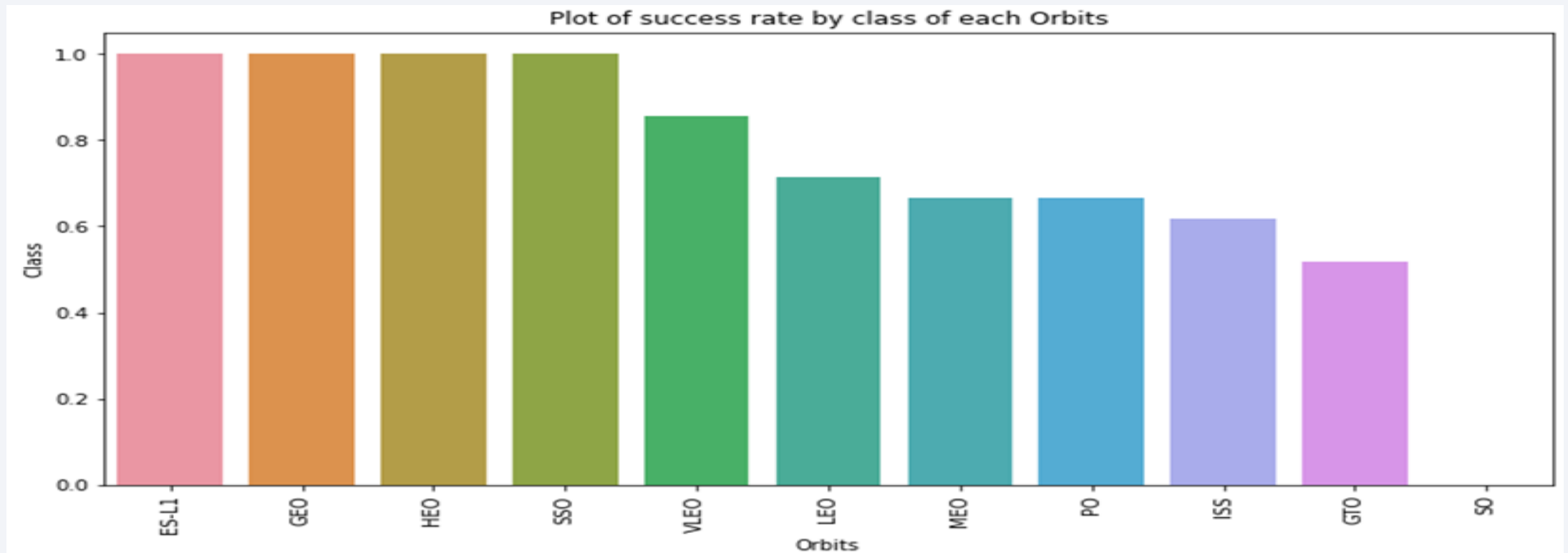
# Payload vs. Launch Site

- The greater the payload mass for launch site CCAFS SLC 40, 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 had the most success rate.



Plot of success rate by class of each Orbits

# 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.



Plot of launch success yearly trend

# All Launch Site Names

- We used the key word DISTINCT to show only unique launch sites from the SpaceX data.

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

```
In [10]:   task_1 = '''
                    SELECT DISTINCT LaunchSite
                    FROM SpaceX
           '''
           create_pandas_df(task_1, database=conn)
```

Out[10]:

|   | launchsite  |
|---|-------------|
| 0 | KSC LC-39A  |
| 1 | CCAFS LC-40 |
| 2 | CCAFS SLC-40 |
| 3 | VAFB SLC-4E |

# Launch Site Names Begin with 'CCA'

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

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

```
In [11]:    task_2 = '''
                SELECT *
                FROM SpaceX
                WHERE LaunchSite LIKE 'CCA%'
                LIMIT 5
                '''
            create_pandas_df(task_2, database=conn)
```

| | date | time | boosterversion | launchsite | payload | payloadmasskg | orbit | customer | missionoutcome | landingoutcome |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 2010-04-06 | 18:45:00 | F9 v1.0 B0003 | CCAFS LC-40 | Dragon Spacecraft Qualification Unit | 0 | LEO | SpaceX | Success | Failure (parachute) |
| 1 | 2010-08-12 | 15:43:00 | F9 v1.0 B0004 | CCAFS LC-40 | Dragon demo flight C1, two CubeSats, barrel of... | 0 | LEO (ISS) | NASA (COTS) NRO | Success | Failure (parachute) |
| 2 | 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 |
| 3 | 2012-08-10 | 00:35:00 | F9 v1.0 B0006 | CCAFS LC-40 | SpaceX CRS-1 | 500 | LEO (ISS) | NASA (CRS) | Success | No attempt |
| 4 | 2013-01-03 | 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)

```
In [12]: task_3 = '''
             SELECT SUM(PayloadMassKG) AS Total_PayloadMass
             FROM SpaceX
             WHERE Customer LIKE 'NASA (CRS)'
             '''
         create_pandas_df(task_3, database=conn)
```

Out[12]:

|   | total_payloadmass |
|---|---|
| 0 | 45596 |

# Average Payload Mass by F9 v1.1

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

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

```
In [13]:   task_4 = '''
                SELECT AVG(PayloadMassKG) AS Avg_PayloadMass
                FROM SpaceX
                WHERE BoosterVersion = 'F9 v1.1'
                '''
           create_pandas_df(task_4, database=conn)
```

```
Out[13]:      avg_payloadmass

           0          2928.4
```

# First Successful Ground Landing Date

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

```
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 |

# Total Number of Successful and Failure Mission Outcomes

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

List the total number of successful and failure mission outcomes

```
In [16]:  task_7a = '''
              SELECT COUNT(MissionOutcome) AS SuccessOutcome
              FROM SpaceX
              WHERE MissionOutcome LIKE 'Success%'
              '''

          task_7b = '''
              SELECT COUNT(MissionOutcome) AS FailureOutcome
              FROM SpaceX
              WHERE MissionOutcome LIKE 'Failure%'
              '''
          print('The total number of successful mission outcome is:')
          display(create_pandas_df(task_7a, database=conn))
          print()
          print('The total number of failed mission outcome is:')
          create_pandas_df(task_7b, database=conn)
```

The total number of successful mission outcome is:

|   | successoutcome |
|---|---|
| 0 | 100 |

The total number of failed mission outcome is:

```
Out[16]:
```

|   | failureoutcome |
|---|---|
| 0 | 1 |

# 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.

List the names of the booster_versions which have carried the maximum payload mass. Use a subquery

```
In [17]: task_8 = '''
            SELECT BoosterVersion, PayloadMassKG
            FROM SpaceX
            WHERE PayloadMassKG = (
                            SELECT MAX(PayloadMassKG)
                            FROM SpaceX
                            )
            ORDER BY BoosterVersion
            '''
         create_pandas_df(task_8, database=conn)
```

Out[17]:
| | boosterversion | payloadmasskg |
|---|---|---|
| 0 | F9 B5 B1048.4 | 15600 |
| 1 | F9 B5 B1048.5 | 15600 |
| 2 | F9 B5 B1049.4 | 15600 |
| 3 | F9 B5 B1049.5 | 15600 |
| 4 | F9 B5 B1049.7 | 15600 |
| 5 | F9 B5 B1051.3 | 15600 |
| 6 | F9 B5 B1051.4 | 15600 |
| 7 | F9 B5 B1051.6 | 15600 |
| 8 | F9 B5 B1056.4 | 15600 |
| 9 | F9 B5 B1058.3 | 15600 |
| 10 | F9 B5 B1060.2 | 15600 |
| 11 | F9 B5 B1060.3 | 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

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 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 2010-03-20.

- 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.

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 |

# Launch Sites Proximities Analysis

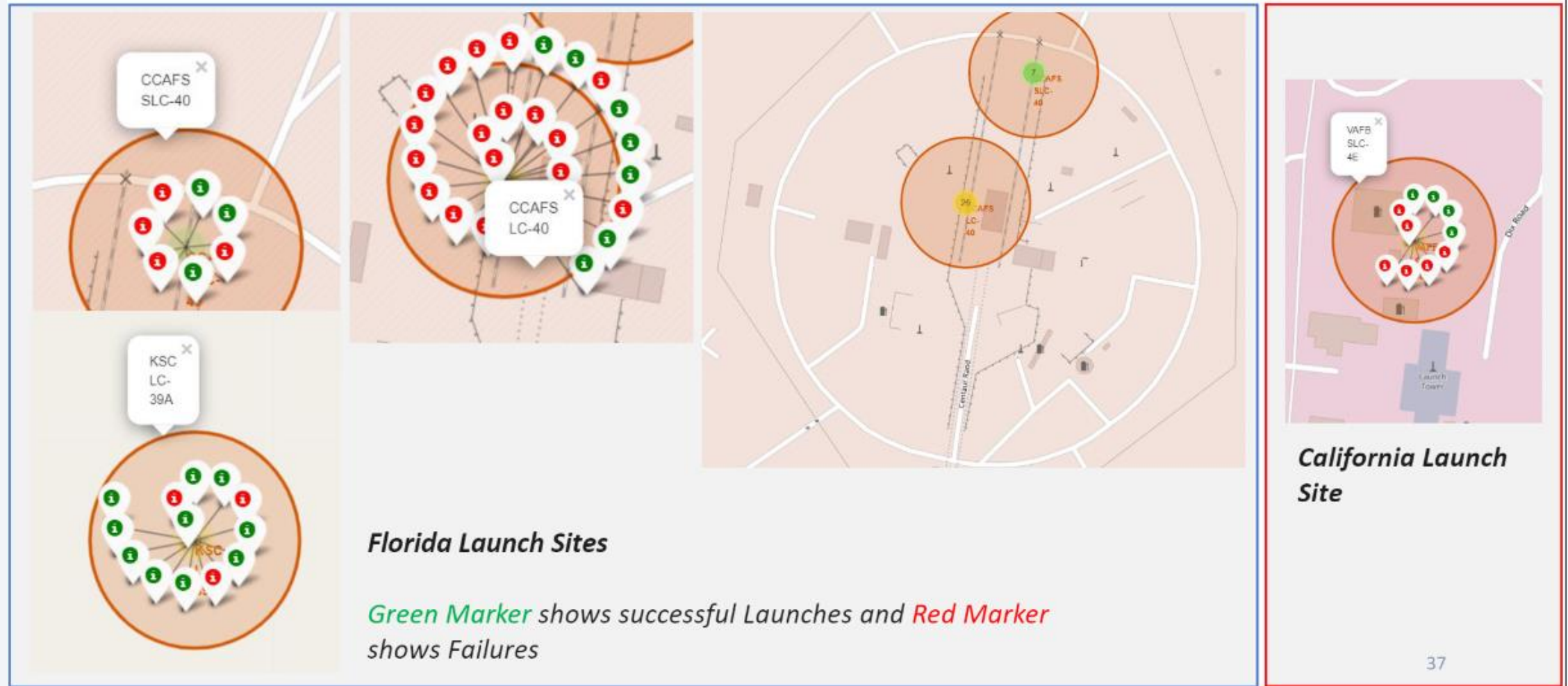All launch sites global map markers

VAFB
SLC-
4E

CCAFS
CCC-
39A

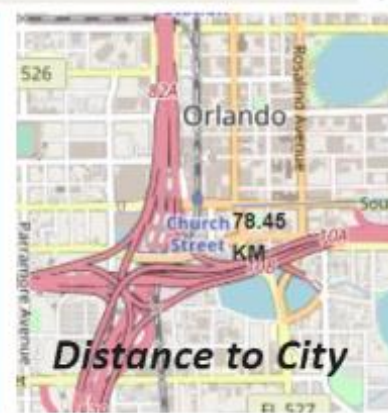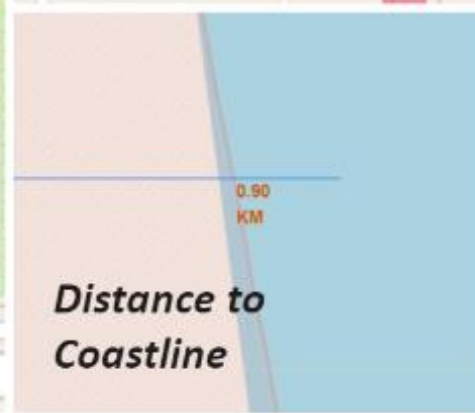We can see that the SpaceX launch sites are in the United States of Florida and California

35

# Markers showing launch sites with color labels



Florida Launch Sites

Green Marker shows successful Launches and Red Marker shows Failures

California Launch Site

37

# Launch Site distance to landmarks



Distance to Railway Station

Distance to closest Highway

Distance to Coastline

Distance to City

Distance to coast

- 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
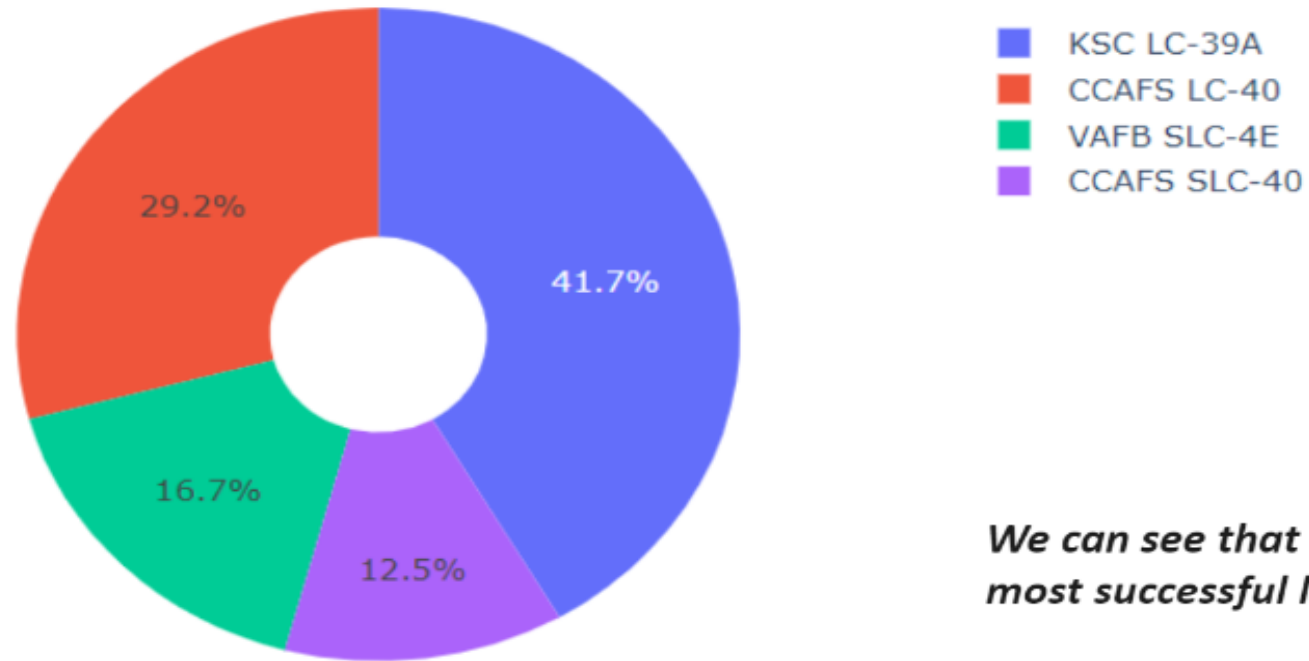- Do launch sites keep certain distance away from cities? Yes

37

Section 4

Build a Dashboard
with Plotly Dash

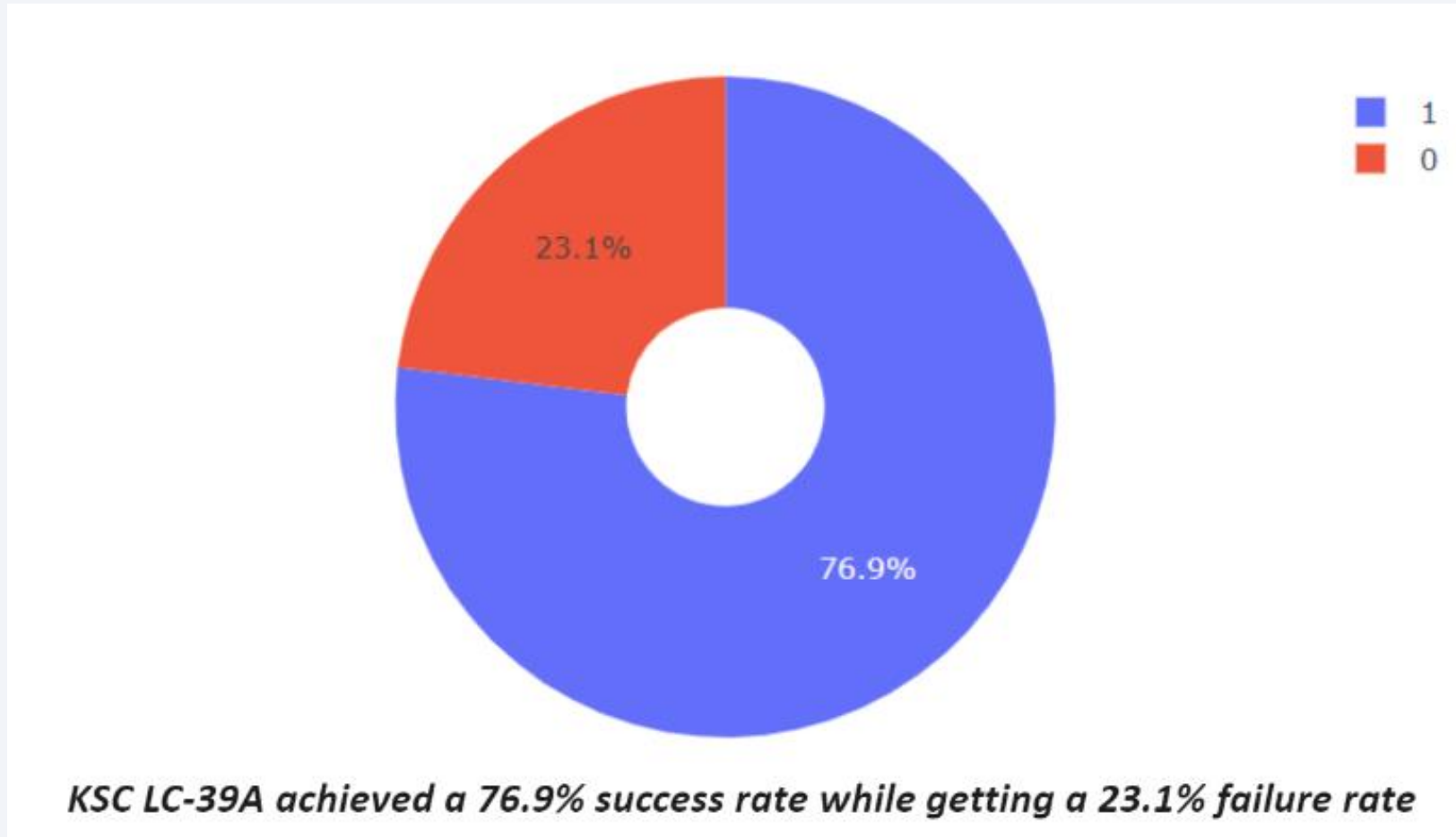# Pie chart showing the success percentage achieved by each launch site



## Total Success Launches By all sites

- 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

*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 ratio



23.1%

76.9%

1

0

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

# Predictive Analysis (Classification)

# Classification Accuracy

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

```python
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_)
```
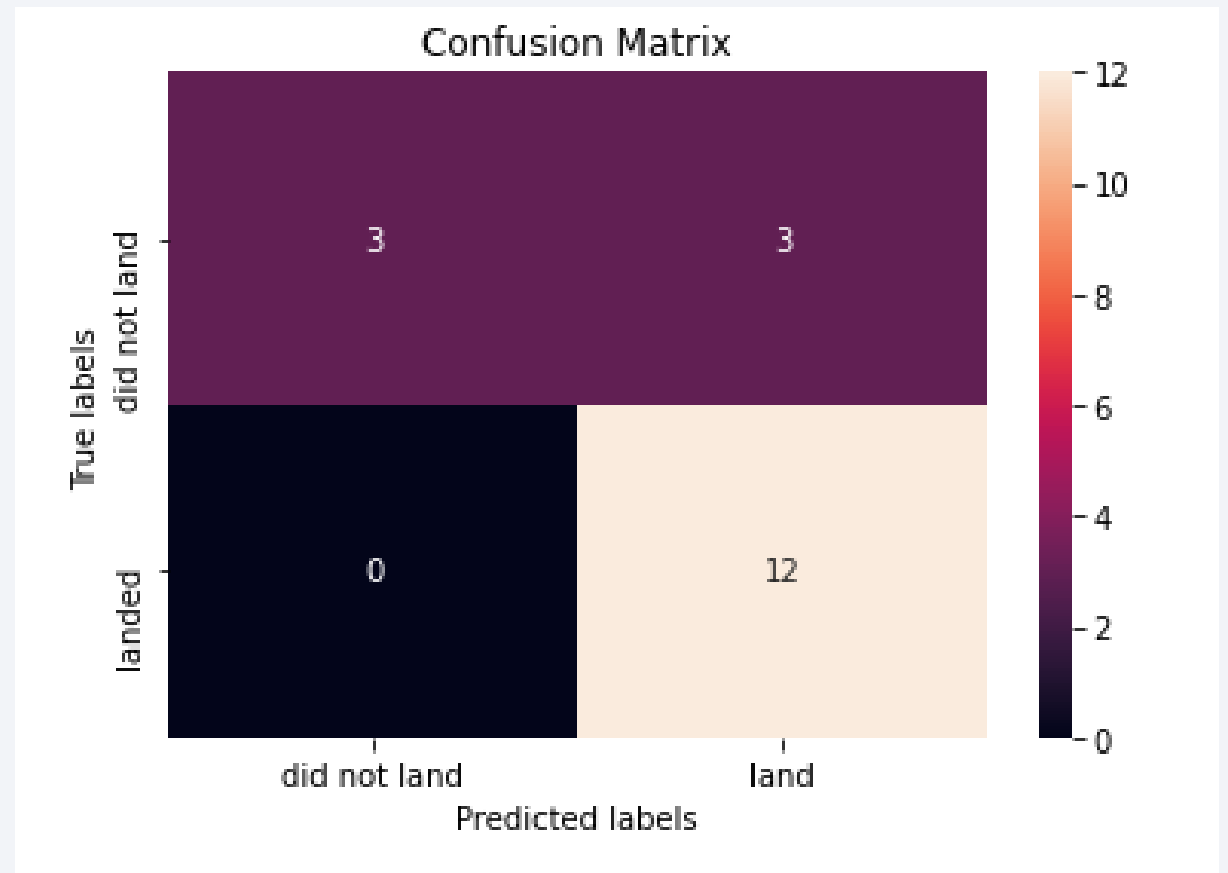
```
Best model is DecisionTree with a score of 0.8732142857142856
Best params is : {'criterion': 'gini', 'max_depth': 6, 'max_features': 'auto', 'min_samples_leaf': 2, 'min_samples_split': 5, 'splitter': 'random'}
```

# 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

- Conclude that the success rate at a launch site increases with the number of flights at that site.

- Observe that the launch success rate has been increasing since 2013 and has continued until 2020.

- Note that the orbits ES-L1, GEO, HEO, SSO, and VLEO had the highest success rates.

- Identify KSC LC-39A as the launch site with the most successful launches.

- Determine that the Decision tree classifier is the best machine learning algorithm for this particular task.

Thank you!