# Outline

- Executive Summary

- Introduction

- Methodology

- Results

- Conclusion

- Appendix

# Executive Summary

**Summary of Methodologies**

The research attempts to identify the factors for a successful rocket landing.
Our methodologies include:
- Collect data using SpaceX REST API and web scraping techniques
- Wrangle data to create success/fail outcome variable
- Explore data with data visualization techniques, considering the following factors:    payload, launch site, flight number and yearly trend
- Analyze the data with SQL, calculating the following statistics: total payload, payload   range for successful launches, and total of successful and failed outcomes
- Explore launch site success rates and proximity to geographical markers
- Visualize the launch sites with the most success and successful payload ranges
- Build Models to predict landing outcomes using logistic regression, support vector
  machine (SVM), decision tree and K-nearest neighbor (KNN)

**Summary of Results**
- EDA:
    - Launch success has improved over time
    - KSC LC-39A has the highest success rate among landing sites
    - Orbits ES-L1, GEO, HEO, and SSO have a 100% success rate
- Geographical Analysis (Folium):
    - Most launch sites are near the equator, and all are close to the coast.
- Predictive Analysis:
    - All models performed similarly on the test set, i.e. accuracy of 83%.

# Introduction

- The advancements in space technology have significantly transformed the commercial space industry, with companies like SpaceX leading the way in cost-effective space
  exploration. One of SpaceX's key innovations is the reusability of the Falcon 9 first-stage booster, which dramatically reduces the cost of space missions. A typical Falcon  9 launch costs $62 million, whereas other space launch providers charge over $165 million per launch. The ability to successfully land and reuse the first-stage booster is a crucial factor in these cost savings.

- This capstone project focuses on predicting the successful landing of the Falcon 9 first stage. By developing a predictive model, we aim to assess the probability of a successful landing and explore its impact on launch costs. This information is particularly valuable for competing aerospace companies seeking to evaluate the feasibility of bidding against SpaceX for commercial satellite launches.
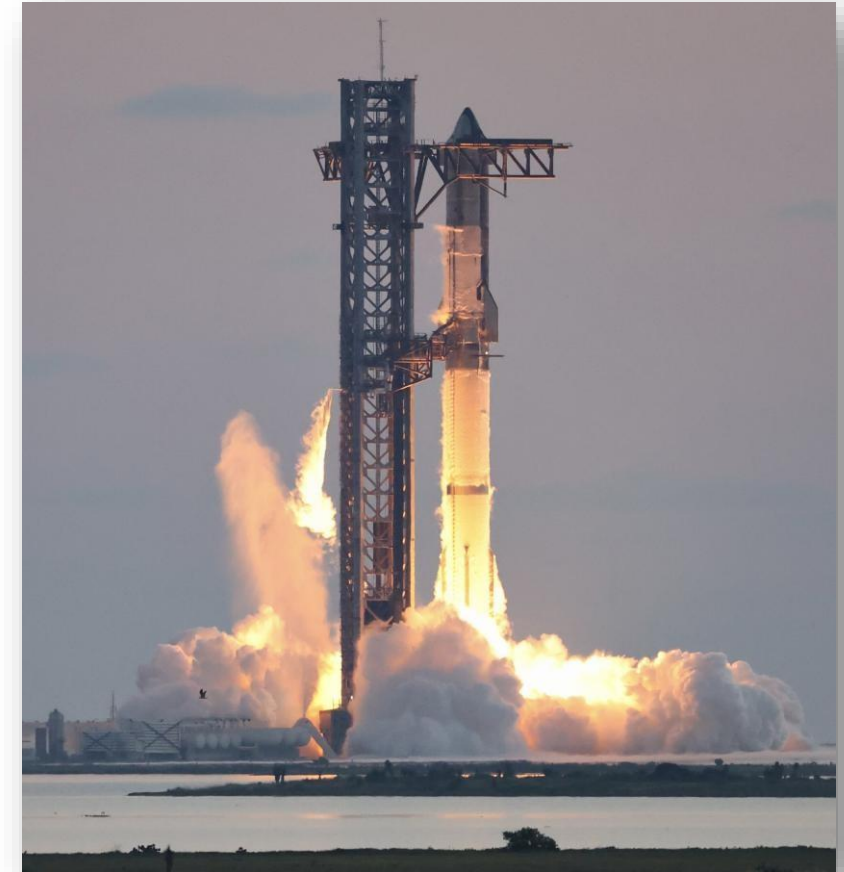
Section 1

# Methodology

# Methodology

## Executive Summary

- Data collection methodology: through SpaceX API and web scrapping Wikipedia.

- Perform data wrangling

- Perform exploratory data analysis (EDA) using visualization and SQL

- Perform interactive visual analytics using Folium and Plotly Dash

- Perform predictive analysis using classification models: Logistic Regression, SVM, KNN, and Decision Tree.
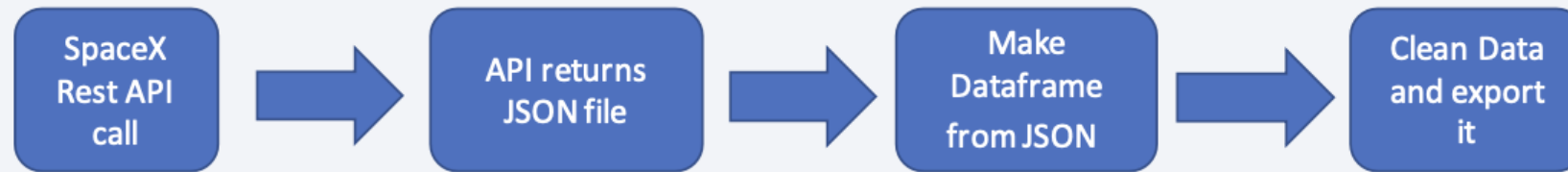
# Data Collection

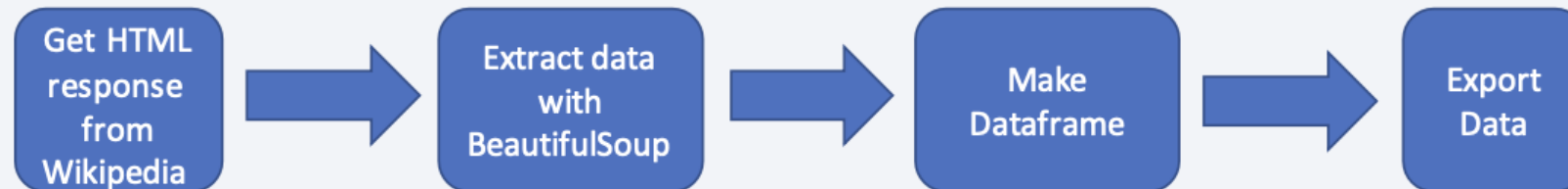- Datasets are collected from Rest SpaceX API and webscrapping Wikipedia

  - The information obtained by the API are rocket, launches, payload information.

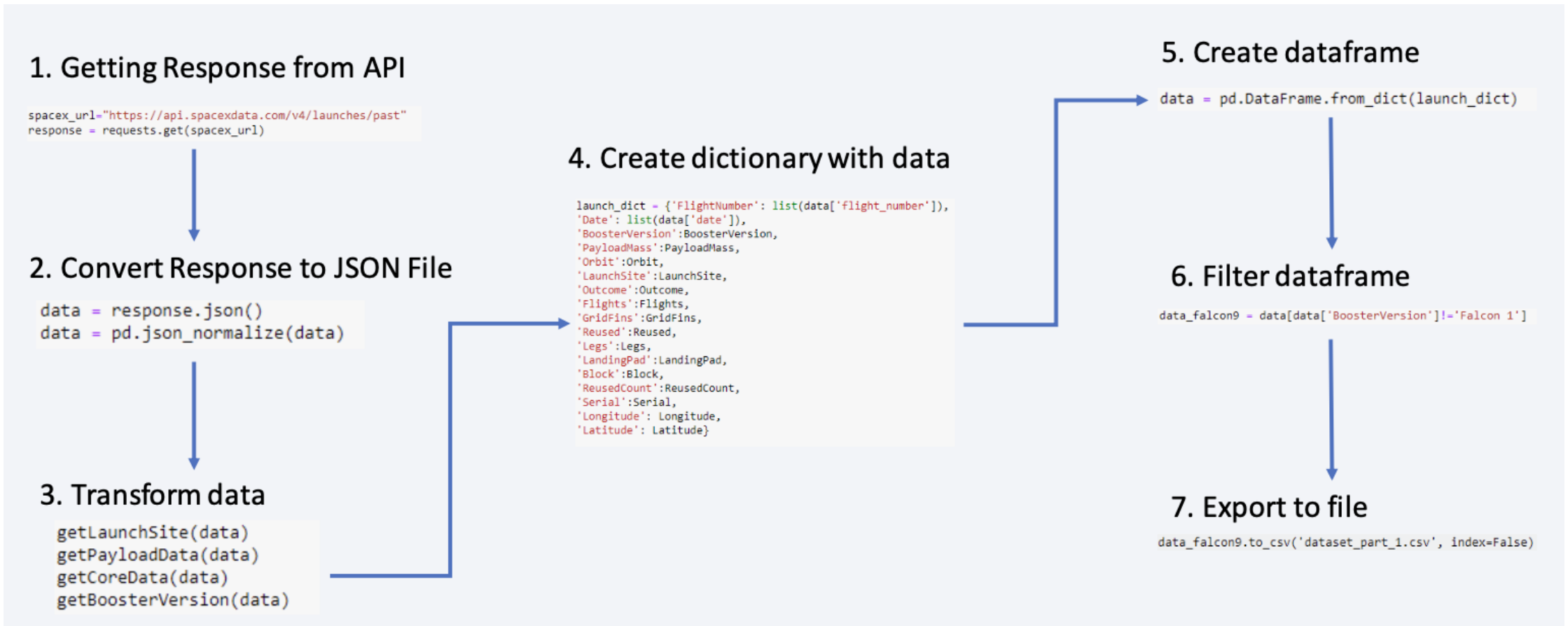    - The Space X REST API URL is api.spacexdata.com/v4/

| SpaceX Rest API call | → | API returns JSON file | → | Make Dataframe from JSON | → | Clean Data and export it |

  - The information obtained by the webscrapping of Wikipedia are launches, landing, payload information.

    - URL is https://en.wikipedia.org/w/index.php?title=List_of_Falcon_9_and_Falcon_Heavy_launches&oldid=1027686922

| Get HTML response from Wikipedia | → | Extract data with BeautifulSoup | → | Make Dataframe | → | Export Data |

# Data Collection – SpaceX API

## 1. Getting Response from API

```
spacex_url="https://api.spacexdata.com/v4/launches/past"
response = requests.get(spacex_url)
```

## 2. Convert Response to JSON File

```
data = response.json()
data = pd.json_normalize(data)
```

## 3. Transform data

```
getLaunchSite(data)
getPayloadData(data)
getCoreData(data)
getBoosterVersion(data)
```

## 4. Create dictionary with data

```
launch_dict = {'FlightNumber': list(data['flight_number']),
'Date': list(data['date']),
'BoosterVersion':BoosterVersion,
'PayloadMass':PayloadMass,
'Orbit':Orbit,
'LaunchSite':LaunchSite,
'Outcome':Outcome,
'Flights':Flights,
'GridFins':GridFins,
'Reused':Reused,
'Legs':Legs,
'LandingPad':LandingPad,
'Block':Block,
'ReusedCount':ReusedCount,
'Serial':Serial,
'Longitude': Longitude,
'Latitude': Latitude}
```

## 5. Create dataframe

```
data = pd.DataFrame.from_dict(launch_dict)
```

## 6. Filter dataframe

```
data_falcon9 = data[data['BoosterVersion']!='Falcon 1']
```

## 7. Export to file

```
data_falcon9.to_csv('dataset_part_1.csv', index=False)
```

The link to the notebook is

https://github.com/snc4shiva/Data-Science-Capstone-Project/blob/main/Lab%20-%20Data%20Collection%20API.ipynb

8

# Data Collection - Scraping

## 1. Getting Response from HTML

```python
response = requests.get(static_url)
```

## 2. Create BeautifulSoup Object

```python
soup = BeautifulSoup(response.text, "html5lib")
```

## 3. Find all tables

```python
html_tables = soup.findAll('table')
```

## 4. Get column names

```python
for th in first_launch_table.find_all('th'):
    name = extract_column_from_header(th)
    if name is not None and len(name) > 0 :
        column_names.append(name)
```

## 5. Create dictionary

```python
launch_dict= dict.fromkeys(column_names)

# Remove an irrelvant column
del launch_dict['Date and time ( )']

# Let's initial the launch_dict with each value to be an empty list
launch_dict['Flight No.'] = []
launch_dict['Launch site'] = []
launch_dict['Payload'] = []
launch_dict['Payload mass'] = []
launch_dict['Orbit'] = []
launch_dict['Customer'] = []
launch_dict['Launch outcome'] = []
# Added some new columns
launch_dict['Version Booster']=[]
launch_dict['Booster landing']=[]
launch_dict['Date']=[]
launch_dict['Time']=[]
```

## 6. Add data to keys

```python
extracted_row = 0
#Extract each table
for table_number,table in enumerate(soup.find_all
    # get table row
    for rows in table.find_all("tr"):
        #check to see if first table heading is a
        if rows.th:
            if rows.th.string:
                flight_number=rows.th.string.stri
                flag=flight_number.isdigit()
```

**See notebook for the rest of code**

## 7. Create dataframe from dictionary

```python
df=pd.DataFrame(launch_dict)
```

## 8. Export to file

```python
df.to_csv('spacex_web_scraped.csv', index=False)
```

The link to the notebook is

https://github.com/snc4shiva/Data-Science-Capstone-Project/blob/main/Lab%20-%20Data%20Collection%20with%20Web%20Scraping.ipynb

# Data Wrangling

- In the dataset, there are several cases where the booster did not land successfully.
  - True Ocean, True RTLS, and True ASDS mean the mission has been successful.
  - False Ocean, False RTLS, False ASDS means the mission was a failure.

- We need to transform string variables into categorical variables where 1 means the mission has been successful and 0 means the mission was a failure.

- The link to the notebook is

https://github.com/snc4shiva/Data-Science-Capstone-Project/blob/main/Lab%20-%20Data%20Wrangling.ipynb



**1. Calculate launches number for each site**
```
df['LaunchSite'].value_counts()

CCAFS SLC 40    55
KSC LC 39A      22
VAFB SLC 4E     13
Name: LaunchSite, dtype: int64
```

**2. Calculate the number and occurence of each orbit**
```
df['Orbit'].value_counts()

GTO     27
ISS     21
VLEO    14
PO       9
LEO      7
SSO      5
MEO      3
SO       1
ES-L1    1
HEO      1
GEO      1
Name: Orbit, dtype: int64
```

**3. Calculate number and occurrence of mission outcome per orbit type**
```
landing_outcomes = df['Outcome'].value_counts()
landing_outcomes

True ASDS     41
None None     19
True RTLS     14
False ASDS     6
True Ocean     5
None ASDS      2
False Ocean    2
False RTLS     1
Name: Outcome, dtype: int64
```

**4. Create landing outcome label from Outcome column**
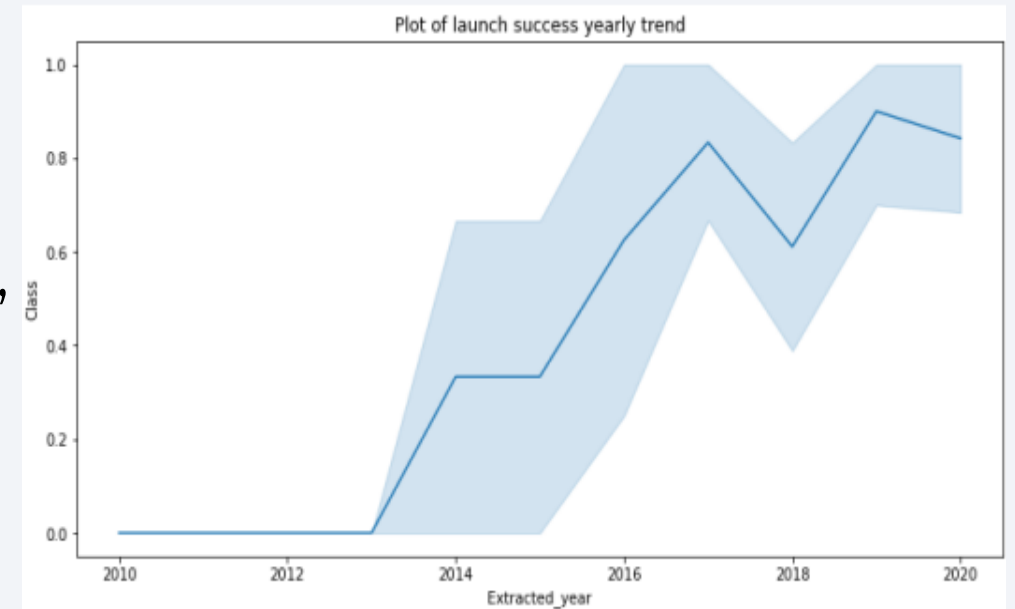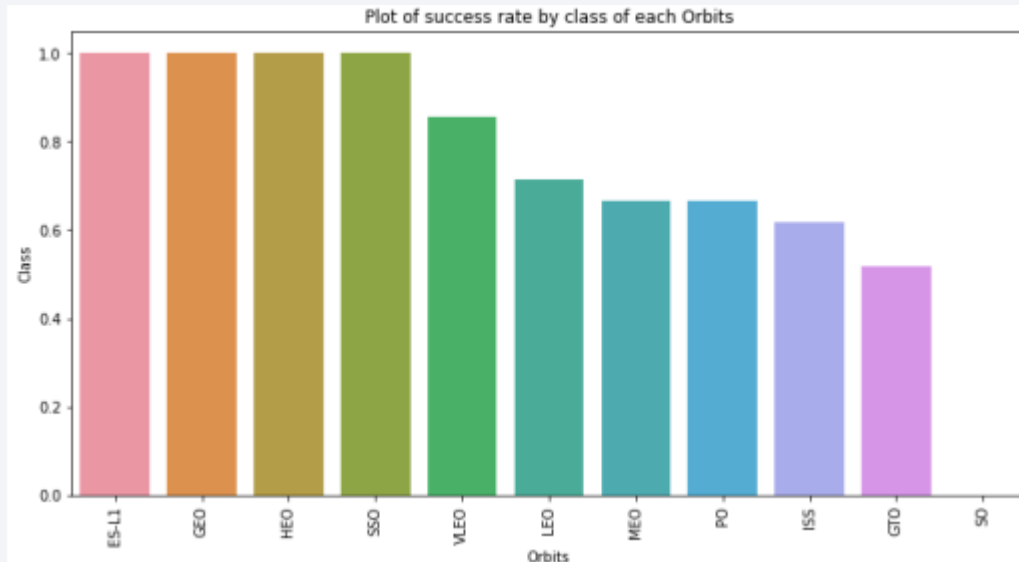```
landing_class = []
for key,value in df["Outcome"].items():
    if value in bad_outcomes:
        landing_class.append(0)
    else:
        landing_class.append(1)
df['Class']=landing_class
```

**5. Export to file**
```
df.to_csv("dataset_part_2.csv", index=False)
```

# EDA with Data Visualization

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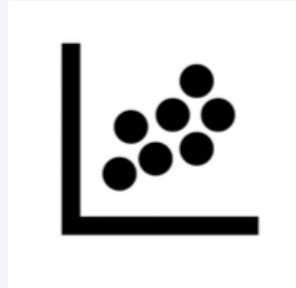
Plot of launch success yearly trend

Plot of success rate by class of each Orbits

- The Link to the notebook is https://github.com/snc4shiva/Data-Science-Capstone-Project/blob/main/Lab%20-%20EDA%20with%20Visualization.ipynb

# EDA with Data Visualization

- **Scatter Graphs**

  - Flight Number vs. Payload Mass

  - Flight Number vs. Launch Site

  - Payload vs. Launch Site

  - Orbit vs. Flight Number

  - Payload vs. Orbit Type

  - Orbit vs. Payload Mass

*Scatter plots show relationship between variables. This relationship is called the correlation.*

- **Bar Graph**
  - Success rate vs. Orbit

*Bar graphs show the relationship between numeric and categoric variables.*

- **Line Graph**
  - Success rate vs. Year

*Line graphs show data variables and their trends. Line graphs can help to show global behavior and make prediction for unseen data.*

# 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 the 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 failed 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/snc4shiva/Data-Science-Capstone-Project/blob/main/Lab%20-%20EDA%20with%20SQL.ipynb

# Build an Interactive Map with Folium

- We marked all launch sites and added map objects such as markers, circles, and 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 rates.

- We calculated the distances between a launch site to its proximities. We answered some questions for instance:

  - Are launch sites near railways, highways, and coastlines.

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

- The link to the notebook is

https://github.com/snc4shiva/Data-Science-Capstone-Project/blob/main/06_SpaceX_Interactive_Visual_Analytics_Folium.ipynb

14

# Build a Dashboard with Plotly Dash

- We built an interactive dashboard with Plotly dash.

- We plotted pie charts showing the total launches by certain sites.

- We plotted a scatter graph showing the relationship between Outcome and Payload Mass (Kg) for the different booster versions.

- The link to the notebook is

https://github.com/snc4shiva/Data-Science-Capstone-Project/blob/main/SpaceX_Interactive_Visual_Analytics_Plotly.py

# Predictive Analysis (Classification)

- We loaded the data using NumPy and pandas, transformed the data, and split our data into training and testing.

- We built different machine-learning models and tuned different hyperparameters using GridSearchCV.

- We used accuracy as the metric for our model and 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/snc4shiva/Data-Science-Capstone-Project/blob/main/08_SpaceX_Predictive_Analytics.ipynb

# Results

- Exploratory data analysis results

- Interactive analytics demo in screenshots

- Predictive analysis results
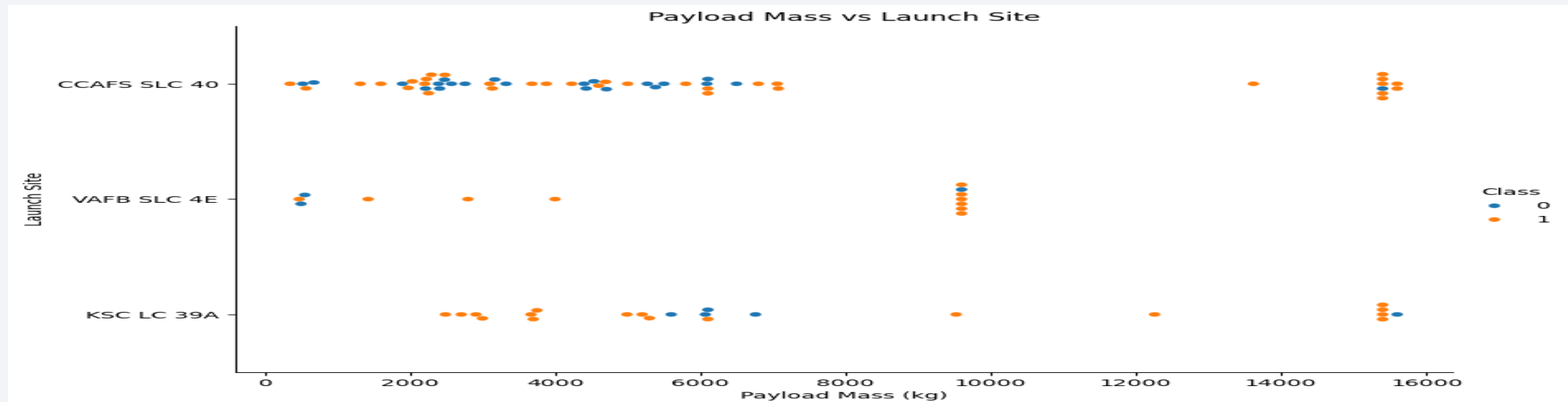
Section 2

# Insights drawn from EDA

# Flight Number vs. Launch Site

- Earlier flights had a lower success rate (blue = fail)

- Later flights had a higher success rate (orange = success)

- Around half of the launches were from the CCAFS SLC 40 launch site

- VAFB SLC 4E and KSC LC 39A have higher success rates

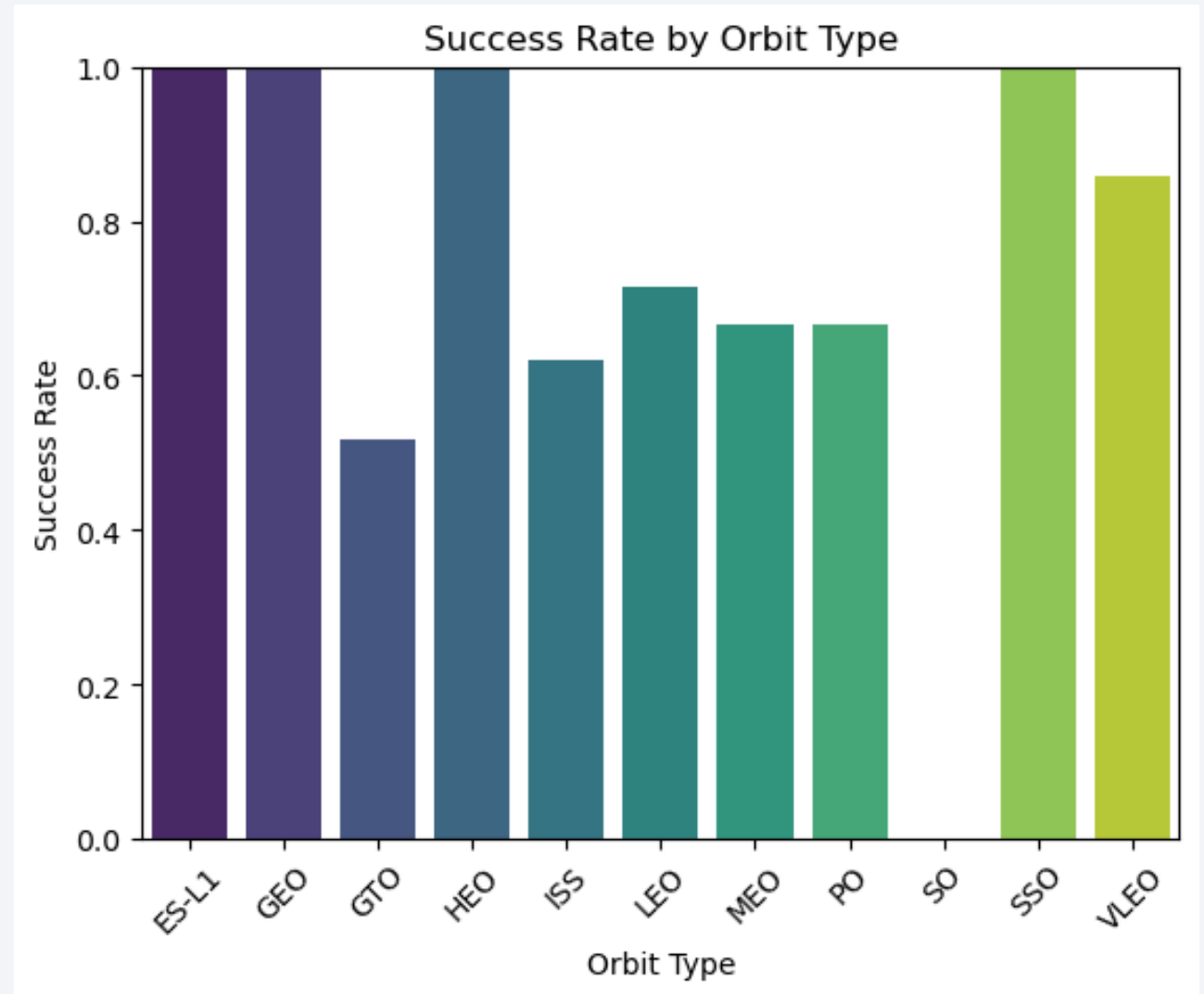- We can infer that new launches have a higher success rate

# Payload vs. Launch Site

- Typically, the higher the payload mass (kg), the higher the success rate
- Most launces with a payload greater than 7,000 kg were successful
- KSC LC 39A has a 100% success rate for launches less than 5,500 kg
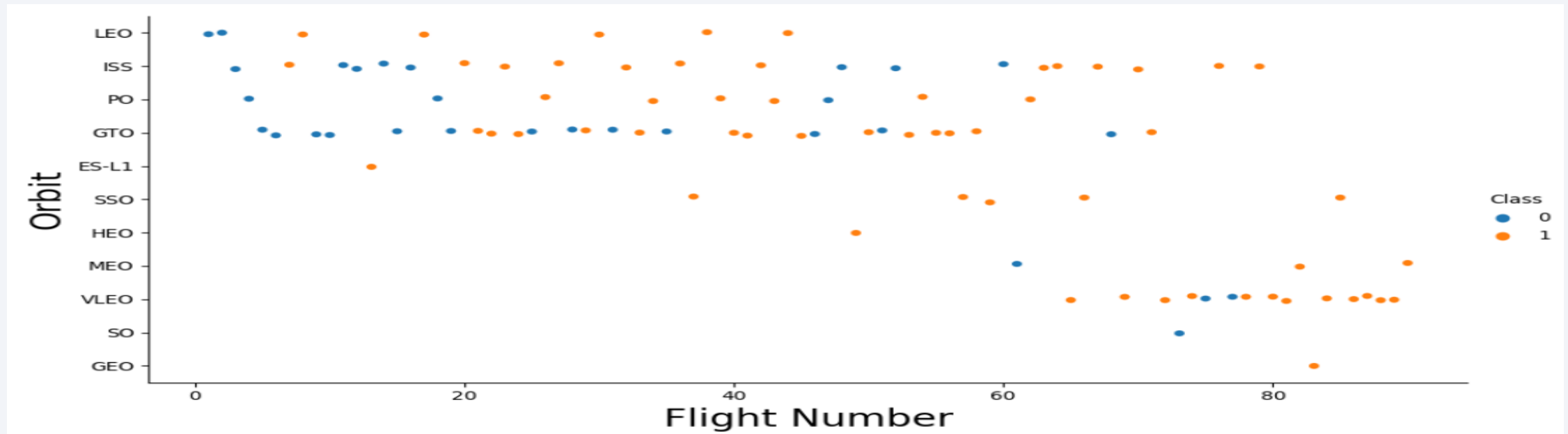- VAFB SKC 4E has not launched anything greater than ~10,000 kg



Payload Mass vs Launch Site

# Success Rate vs. Orbit Type

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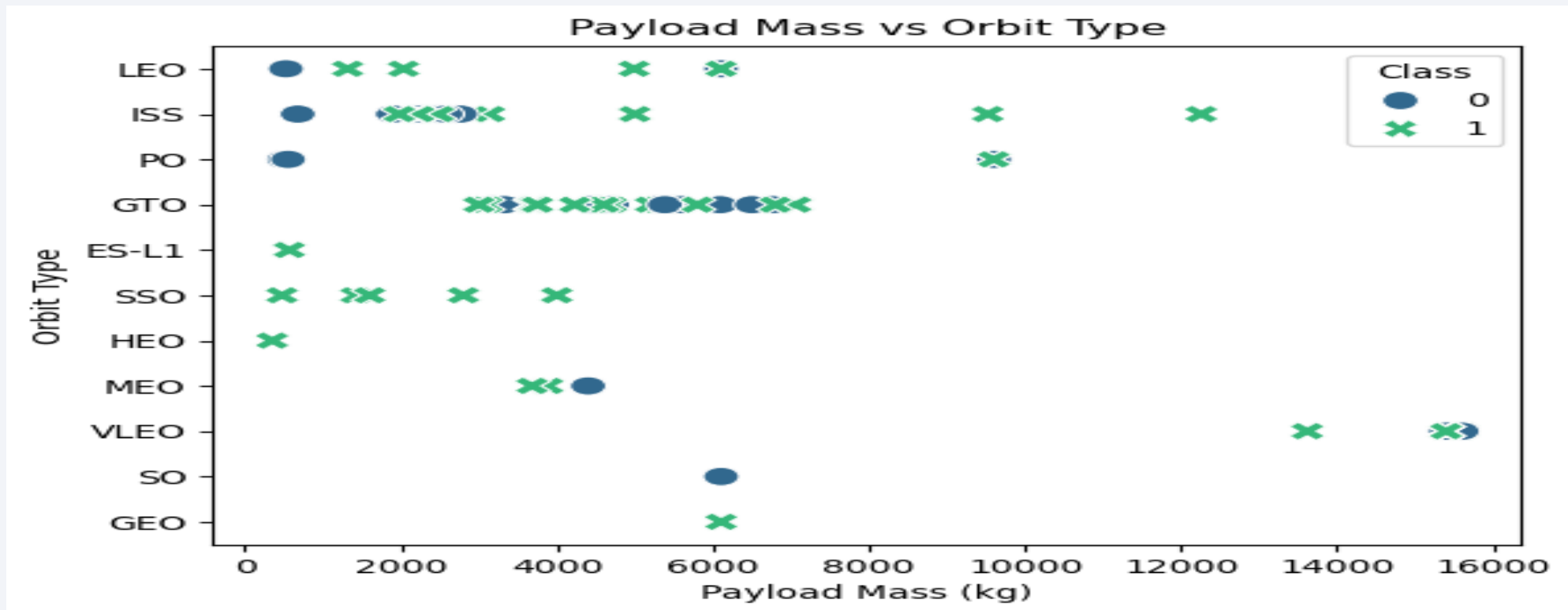


Success Rate by Orbit Type

# 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.
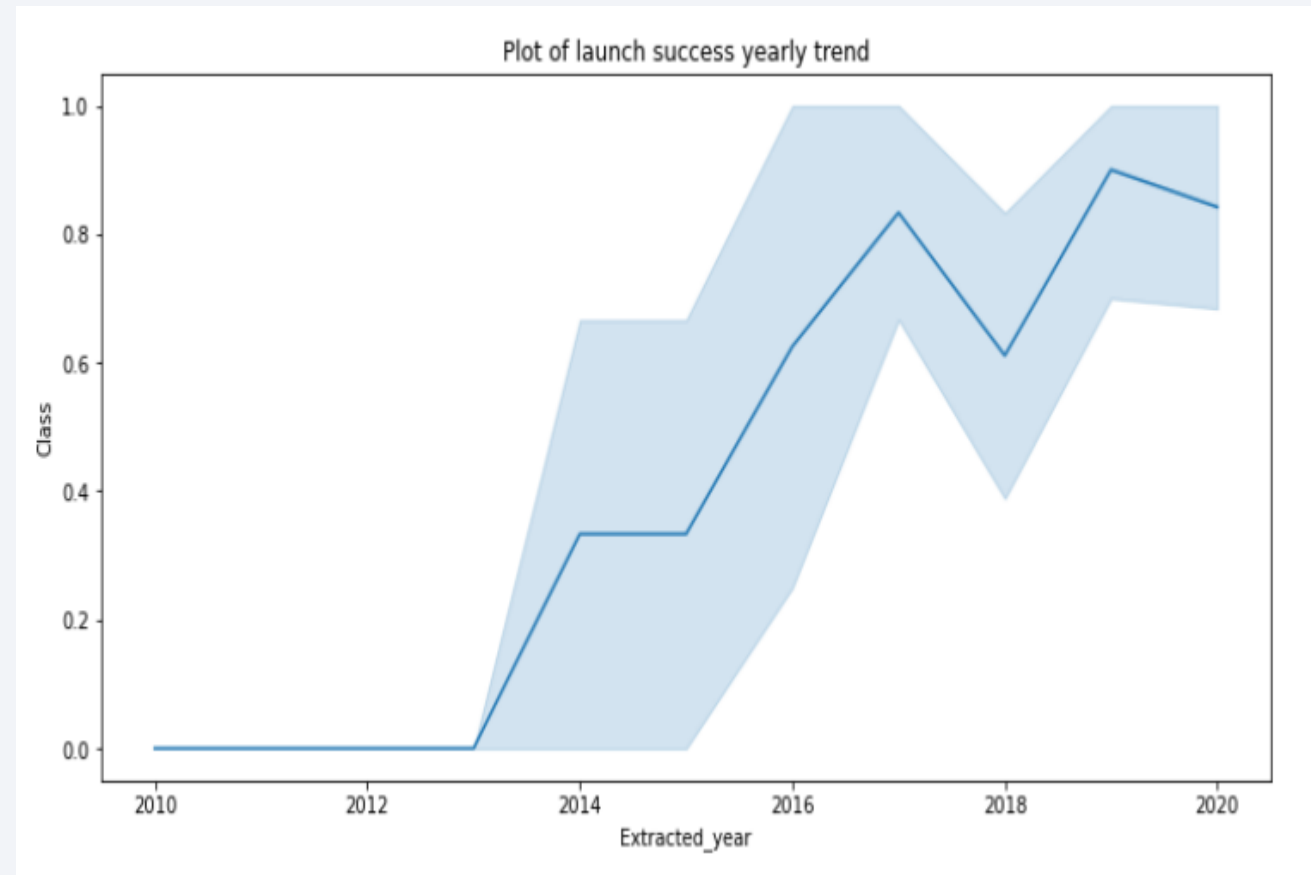
# Payload vs. Orbit Type

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

# Launch Success Yearly Trend

- The success rate improved from 2013-2017 and 2018-2019

- The success rate decreased from 2017-2018 and 2019-2020

- Overall, the success rate has improved since 2013



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.

- The table below shows the 5 records from the launch site CCAFS LC-40.

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)
```

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

| 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 |

# 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

```
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 that have successfully landed on a drone ship and applied the AND condition to determine successful landing with payload mass greater than 4000 but less than 6000

# Total Number of Successful and Failure Mission Outcomes

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 |

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

31

# Boosters Carried Maximum Payload

- We determined the booster that has 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

```
n [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)
```

ut[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 combination of the WHERE clause, LIKE, AND, and BETWEEN conditions to filter for failed landing outcomes in drone ships, their booster versions, and launch site names for the year 2015

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

```
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)
```

| | 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 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 outcomes in descending order.
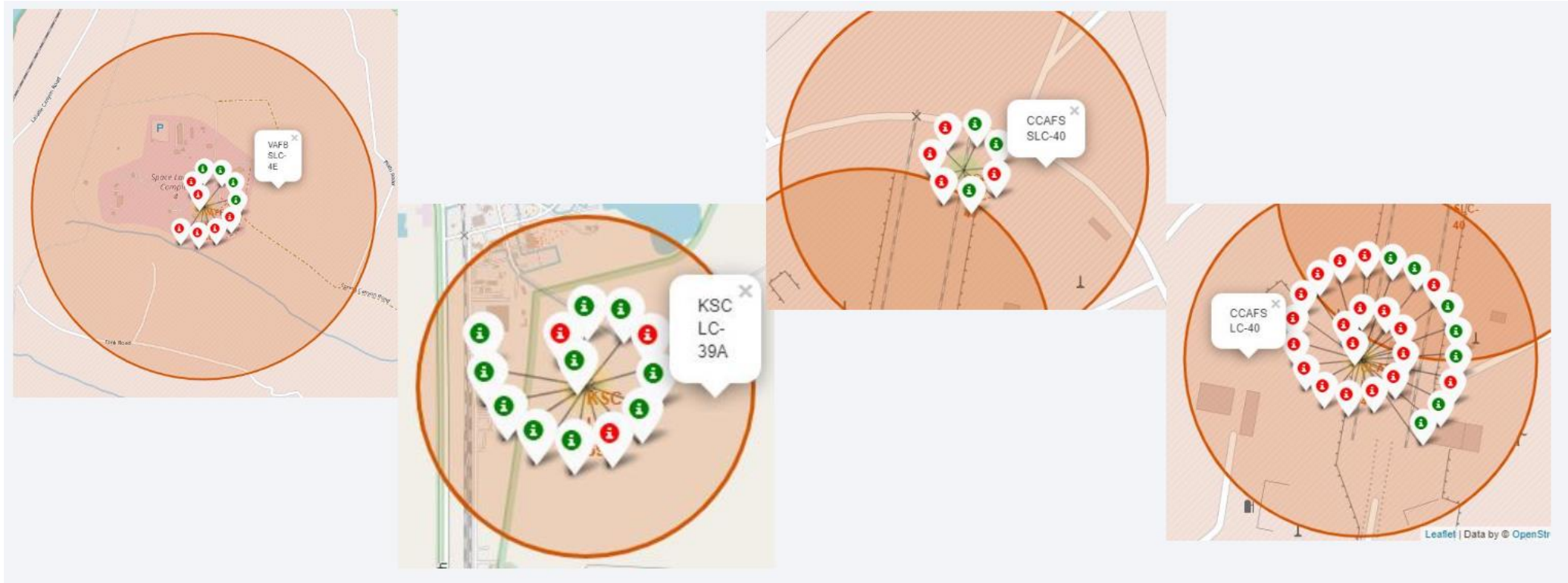
Section 3

# Launch Sites Proximities Analysis

# Folium map – Ground stations



We see that Space X launch sites are located on the coast of the United States

# Folium map – Color Labeled Markers



The green marker represents successful launches. The Red marker represents unsuccessful launches.
We note that KSC LC-39A has a higher launch success rate.

37

# Folium Map – Distances between CCAFS SLC-40 and its proximities



Is CCAFS SLC-40 in close proximity to railways ? Yes
Is CCAFS SLC-40 in close proximity to highways ? Yes
Is CCAFS SLC-40 in close proximity to coastline ? Yes
Do CCAFS SLC-40 keeps certain distance away from cities ? No

Section 4

# Build a Dashboard with Plotly Dash

# Total success by Site



Total Successful Launches by Site

Legend:
- KSC LC-39A
- CCAFS LC-40
- VAFB SLC-4E
- CCAFS SLC-40

Pie chart values:
- 41.7%
- 29.2%
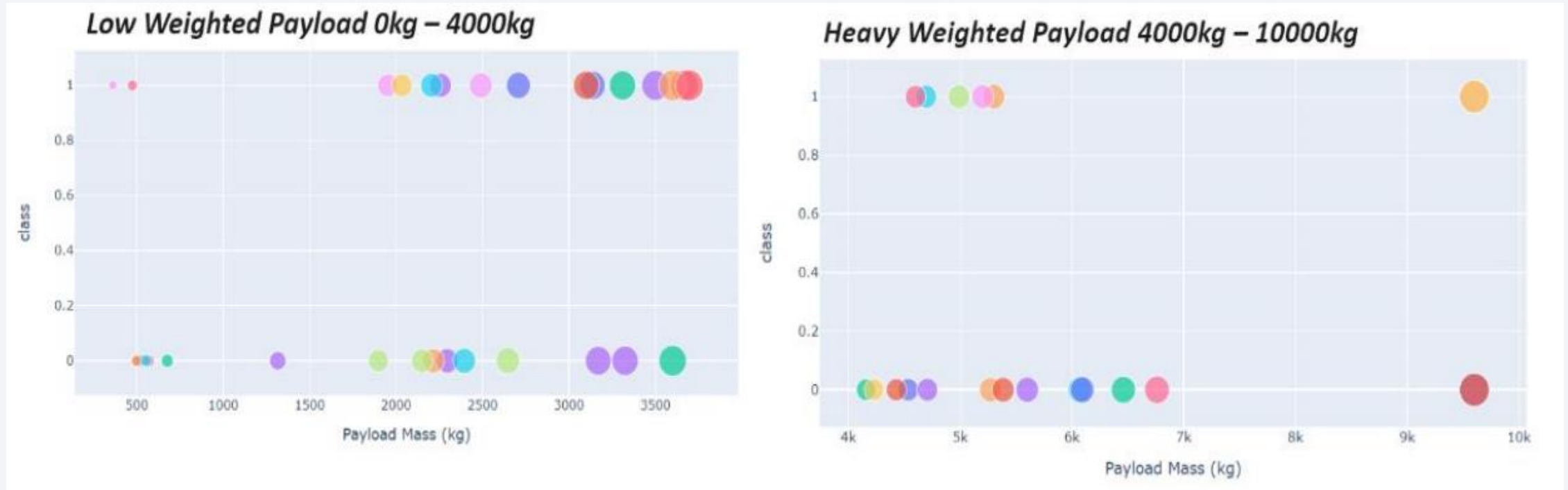- 16.7%
- 12.5%

We see that KSC LC-39A has the best success rate of launches.

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

Total Success/Failure for Site KSC LC-39A



■ 1
■ 0

23.1%

76.9%

We see that KSC LC-39A has 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 success rate for low-weighted payloads is higher than the heavy-weighted payloads.

Section 5

# Predictive Analysis (Classification)

# Classification Accuracy

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

```
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_)
```
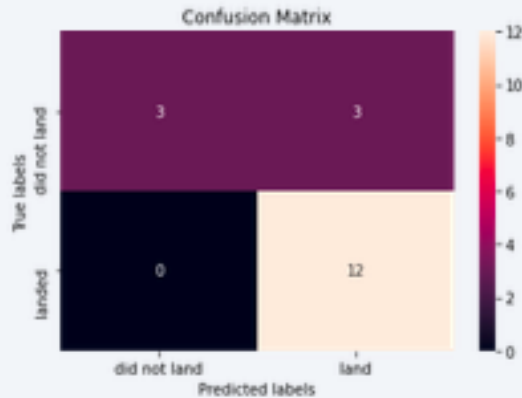
```
Best model is DecisionTree with a score of 0.8732142857142856
Best params is : {'criterion': 'gini', 'max_depth': 6, 'max_features': 'auto', 'min_samples_leaf': 2, 'min_samples_split': 5, 'splitter': 'random'}
```
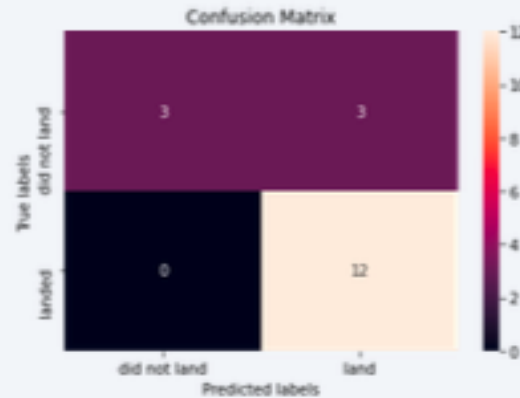
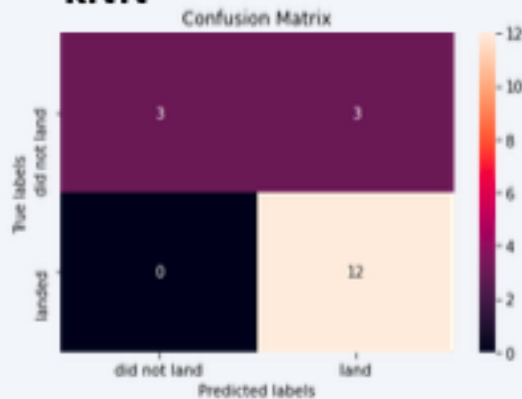|   | Model | Train Accuracy Score | Test Accuracy Score |
|---|-------|---------------------|---------------------|
| 0 | Logistic Regression | 0.846429 | 0.833333 |
| 1 | SVM | 0.848214 | 0.833333 |
| 2 | Decision Tree | 0.875000 | 0.833333 |
| 3 | KNN | 0.848214 | 0.833333 |

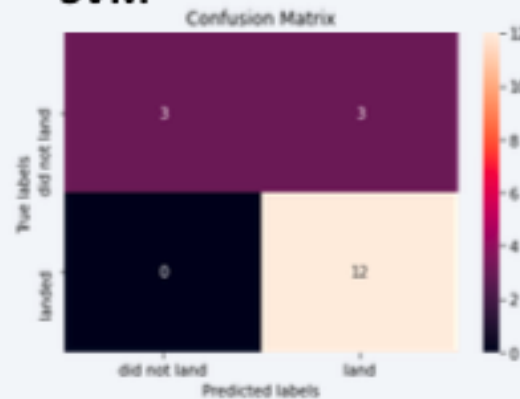# Confusion Matrix



As the test accuracy are all equal, the confusion matrices are also identical. The main problem of these models are false positives.

# Conclusions

- Any of these models can be used for prediction, but simpler models like Logistic Regression or SVM might be preferred due to their interpretability and robustness.

- Decision Tree has the highest training accuracy but does not generalize better than the others.

- If needed, further hyperparameter tuning or trying ensemble methods (like Random Forest) might improve performance



PERFECTING PROPULSIVE LANDING

Thank you!