# Outline

- Executive Summary

- Introduction

- Methodology

- Results

- Conclusion

- Appendix

# Executive Summary

- Summary of methodologies

    - Data Collection via API, Web Scraping

    - Exploratory Data Analysis (EDA) with Data Visualization

    - EDA with SQL

    - Interactive Map with Folium

    - Dashboards with Plotly Dash

    - Predictive Analysis

- Summary of all results

    - Exploratory Data Analysis results

    - Interactive maps and dashboard

    - Predictive results

# Introduction

- SpaceX is a revolutionary company which disrupted the space industry by offering a rocket launches specifically Falcon 9 as low as 62 million dollars; while other providers cost upward of 165 million dollar each. Most of this saving thanks to SpaceX astounding idea to reuse the first stage of the launch by re-land the rocket to be used on the next mission.

- The aim of this project is to predict if the Falcon 9 first stage will successfully land. By determining if the stage will land, we can determine the cost of a launch.

Section 1

# Methodology

# Methodology

Executive Summary

- Data collection methodology:

  - Data was collected using SpaceX REST API and web scrapping from Wikipedia

- Perform data wrangling

  - Describe was using one-hot encoding for 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

- Datasets are collected from Rest SpaceX API and web scrapping through Wikipedia.

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

SpaceX Rest API call → API returns JSON file → Make Data frame From JSON → Clean Data and export it

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

Get HTML Response from Wikipedia → Extract data with BeautifulSoup → Make Data frame → 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)
```

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)
```

# Data Wrangling

- In the dataset, there are several cases where the booster did not land successfully.
  - True Ocean, True RTLS, True ASDS means 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.

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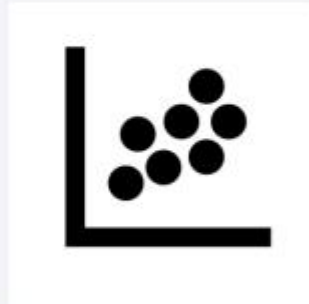
**5. Export to file**
```
df.to_csv("dataset_part_2.csv", index=False)
```

# EDA with Data Visualization

- 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 performed SQL queries to gather and understand data from dataset:

  - Displaying the names of the unique launch sites in the space mission.

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

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

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

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

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

  - List the total number of successful and failure mission outcomes.

  - List the names of the booster versions which have carried the maximum payload mass.

  - List the records which will display the month names, failure landing_outcomes in drone ship, booster versions, launch site for the months in year 2015.

  - Rank the count of successful landiing_outcomes between the date 04-06-2010 and 20-03-2017 in descending order.
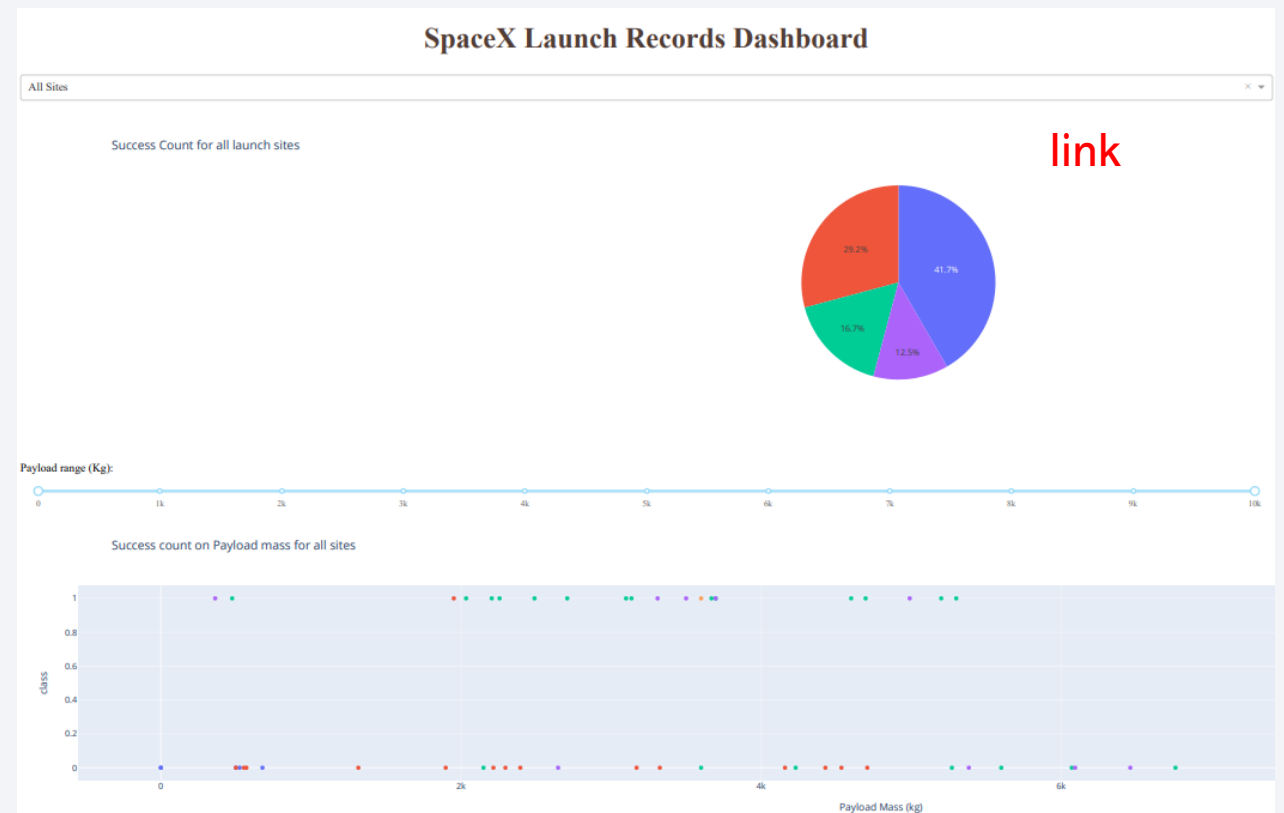
# Build an Interactive Map with Folium

- Folium map object is a map centered on NASA Johnson Space Center at Houson, Texas

  - Red circle at NASA Johnson Space Center's coordinate with label showing its name *(folium.Circle, folium.map.Marker)*.

  - Red circles at each launch site coordinates with label showing launch site name *(folium.Circle, folium.map.Marker, folium.features.DivIcon)*.

  - The grouping of points in a cluster to display multiple and different information for the same coordinates *(folium.plugins.MarkerCluster)*.

  - Markers to show successful and unsuccessful landings. Green for successful landing and Red for unsuccessful landing. *(folium.map.Marker, folium.Icon)*.

  - Markers to show distance between launch site to key locations (railway, highway, coastway, city) and plot a line between them. *(folium.map.Marker, folium.PolyLine, folium.features.DivIcon)*

- These objects are created in order to understand better the problem and the data. We can show easily all launch sites, their surroundings and the number of successful and unsuccessful landings.

link

13

# Build a Dashboard with Plotly Dash

- Dashboard has dropdown, pie chart, rangeslider and scatter plot components
- Dropdown allows a user to choose the launch site or all launch sites(dash_core_components.Dropdown).
- Pie chart shows the total success and the total failure for the launch site chosen with the dropdown component (plotly.express.pie).
- Rangeslider allows a user to select a payload mass in a fixed range (dash_core_components.RangeSlider).
- Scatter chart shows the relationship between two variables, in particular Success vs Payload Mass (plotly.express.scatter).



link

# Predictive Analysis (Classification)

## Building the Model

- Load the dataset into NumPy and Pandas
- Transform the data and then split into training and test datasets
- Decide which type of ML to use
- set the parameters and algorithms to GridSearchCV and fit it to dataset.

## Evaluating the Model

- Check the accuracy for each model
- Get tuned hyperparameters for each type of algorithms.
- plot the confusion matrix.

## Improving the Model

- Use Feature Engineering and Algorithm Tuning

## Find the Best Model

- The model with the best accuracy score will be the best performing model.

link

# Results

- The results will be categorized to 3 main results which is:

    - Exploratory data analysis results

    - Interactive analytics demo in screenshots

    - Predictive analysis results

Section 2

# Insights drawn from EDA
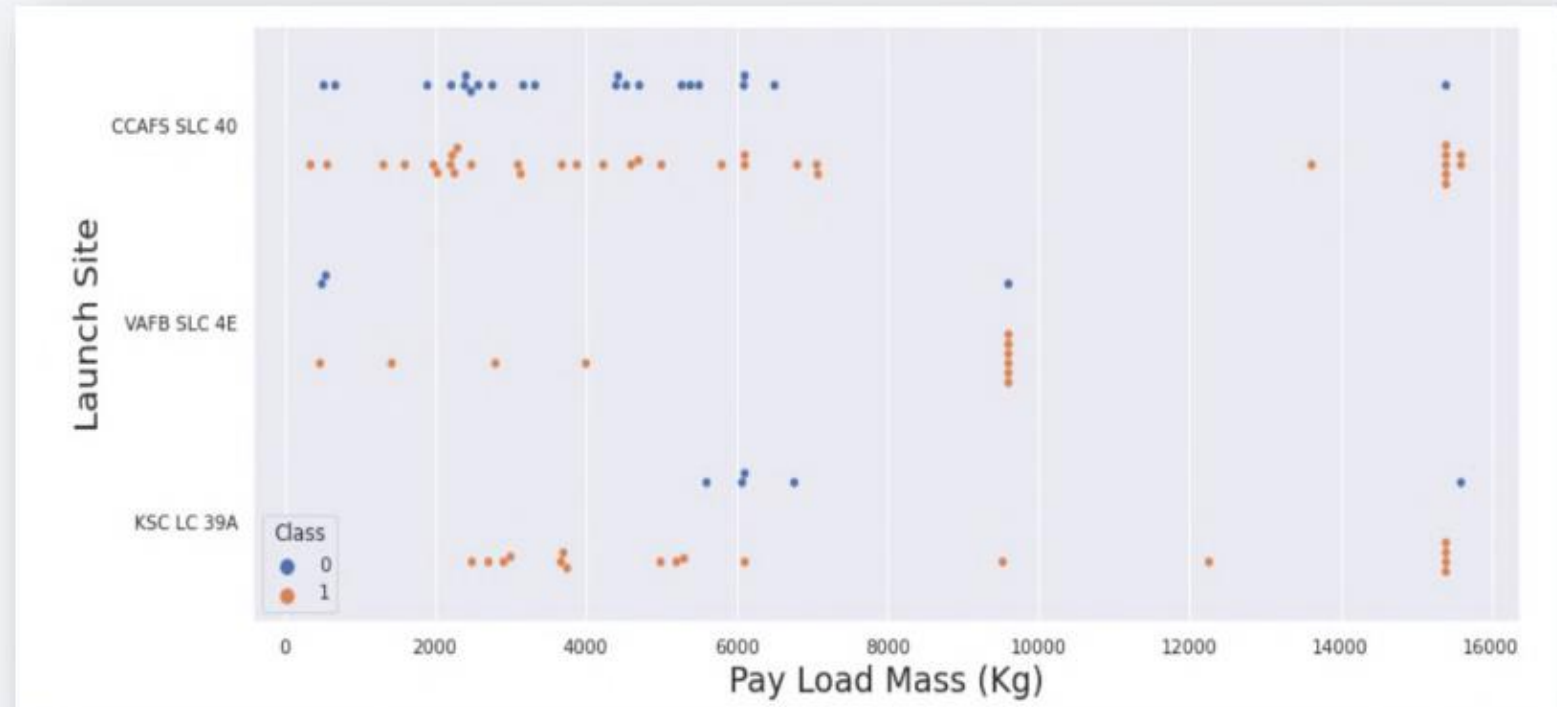
# Flight Number vs. Launch Site



This scatter plot shows that the larger the flights amount of the launch site, the greater the the success rate will be.

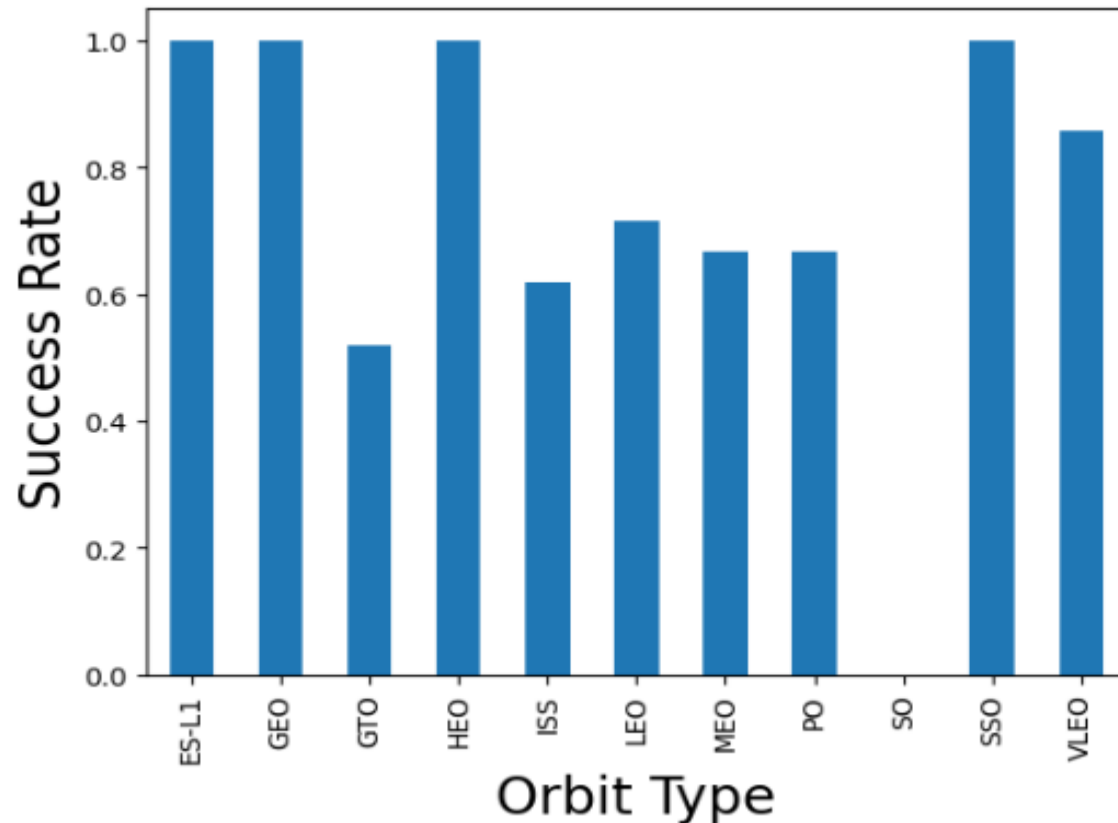However, site CCAFS SLC40 shows the least pattern of this.

# Payload vs. Launch Site

This scatter plot shows once the pay load mass is greater than 7000kg, the probability of the success rate will be highly increased.

However, there is no clear pattern to say the launch site is dependent to the pay load mass for the success rate.
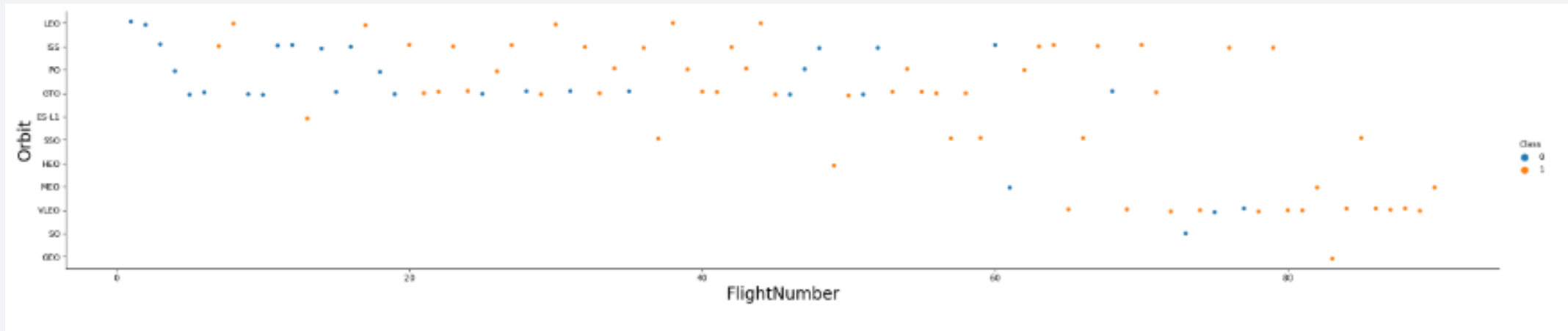
# Success Rate vs. Orbit Type



This figure depicted the possibility of the orbits to influences the landing outcomes as some orbits has 100% success rate such as SSO, HEO, GEO AND ES-L1 while SO orbit produced 0% rate of success.
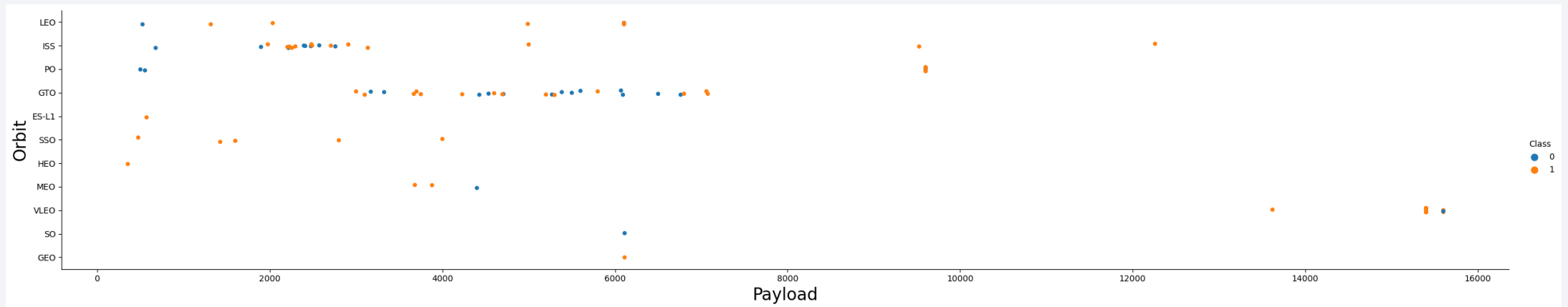
However, deeper analysis show that some of this orbits has only 1 occurrence such as GEO, SO, HEO and ES-L1 which mean this data need more dataset to see pattern or trend before we draw any conclusion.
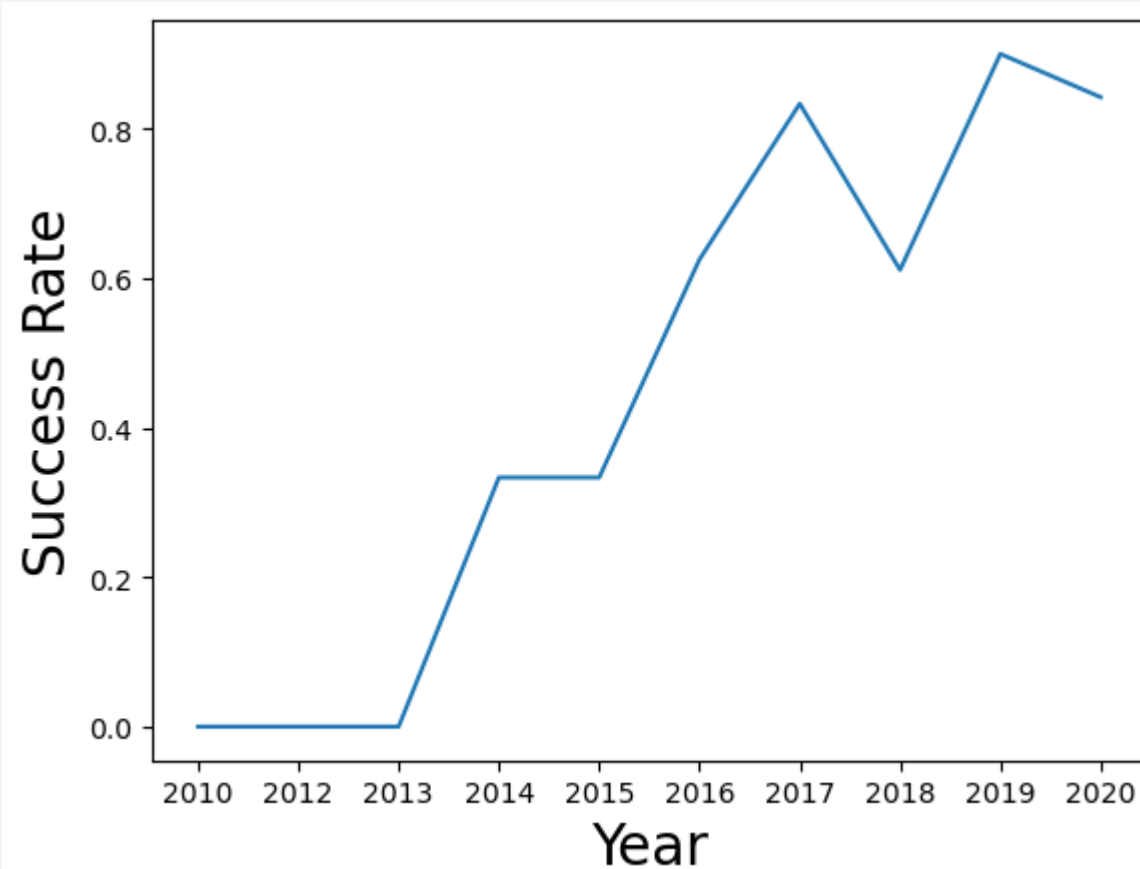
# Flight Number vs. Orbit Type



This scatter plot shows that generally, the larger the flight number on each orbits, the greater the success rate (especially LEO orbit) except for GTO orbit which depicts no relationship between both attributes. Orbit that only has 1 occurrence should also be excluded from above statement as it's needed more dataset.

# Payload vs. Orbit Type



Heavier payload has positive impact on LEO, ISS and P0 orbit. However, it has negative impact on MEO and VLEO orbit. GTO orbit seem to depict no relation between the attributes. Meanwhile, again, SO, GEO and HEO orbit need more dataset to see any pattern or trend.

# Launch Success Yearly Trend



This figures clearly depicted and increasing trend from the year 2013 until 2020. If this trend continue for the next year onward. The success rate will steadily increase until reaching 1/100% success rate.

# All Launch Site Names

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

```
In [5]: %sql SELECT DISTINCT LAUNCH_SITE as "Launch_Sites" FROM SPACEX;

         * ibm_db_sa://zpw86771:***@fbd88901-ebdb-4a4f-a32e-9822b9fb237b.clogj3
         sd0tgtu01qde00.databases.appdomain.cloud:32731/bludb
         Done.

Out[5]:  Launch_Sites

         CCAFS LC-40

         CCAFS SLC-40

         KSC LC-39A

         VAFB SLC-4E
```

# Launch Site Names Begin with 'CCA'

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

```sql
%%sql
SELECT *
FROM SPACEXTBL
WHERE LAUNCH_SITE like 'CCA%'
LIMIT 5;
```

 * db2://vqw61101:***@b0aebb68-94fa-46ec-a1fc-1c999edb6187.c3n41cmd0nqnrk39u98g.databases.appdomain.cloud:31249/bludb
Done.

2]:

| 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

- Total payload carried by boosters from NASA

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

```
In [73]:   %%sql
           SELECT SUM(PAYLOAD_MASS__KG_)
           FROM SPACEXTBL
           WHERE Customer = 'NASA (CRS)';
```

 * db2://vqw61101:***@b0aebb68-94fa-46ec-a1fc-1c999edb6187.c3n41cmd0nqnrk39u98g.databases.appdomain.cloud:31249/bludb
Done.

```
Out[73]:        1

           45596
```
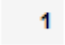
# Average Payload Mass by F9 v1.1

- Average payload mass carried by booster version F9 v1.1

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

```
In [74]:    %%sql
            SELECT AVG(PAYLOAD_MASS__KG_)
            FROM SPACEXTBL
            WHERE Booster_Version LIKE 'F9 v1.0%';

             * db2://vqw61101:***@b0aebb68-94fa-46ec-a1fc-1c999edb6187.c3n41cmd0nqnrk39u98g.databases.appdomain.cloud:31249/bludb
            Done.

Out[74]:      1

            340
```

# First Successful Ground Landing Date

- Dates of the first successful landing outcome on ground pad

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

*Hint:Use min function*

```
In [80]:    %%sql
            SELECT MIN(a.Date)
            FROM (select * from SPACEXTBL) as a
            WHERE a.Landing__Outcome = 'Success (ground pad)';
```

 * db2://vqw61101:***@b0aebb68-94fa-46ec-a1fc-1c999edb6187.c3n41cmd0nqnrk39u98g.databases.appdomain.cloud:31249/bludb
Done.

Out[80]:

| 1 |
|---|
| 2015-12-22 |

# Successful Drone Ship Landing with Payload between 4000 and 6000

- Names of boosters which have successfully landed on drone ship and had payload mass greater than 4000 but less than 6000

```
In [82]:   %%sql
           SELECT BOOSTER_VERSION
           FROM SPACEXTBL
           WHERE Landing__Outcome = 'Success (drone ship)'
               AND 4000 < PAYLOAD_MASS__KG_ < 6000;

            * db2://vqw61101:***@b0aebb68-94fa-46ec-a1fc-1c999edb6187.c3n41cmd0nqnrk39u98g.databases.appdomain.cloud:31249/bludb
           Done.
```

Out[82]:

| booster_version |
|---|
| F9 FT B1021.1 |
| F9 FT B1023.1 |
| F9 FT B1029.2 |
| F9 FT B1038.1 |
| F9 B4 B1042.1 |
| F9 B4 B1045.1 |
| F9 B5 B1046.1 |

# Total Number of Successful and Failure Mission Outcomes

- Total number of successful and failure mission outcomes

*List the total number of successful and failure mission outcomes*

```
In [83]:   %%sql
           SELECT MISSION_OUTCOME, COUNT(MISSION_OUTCOME) AS TOTAL_NUMBER
           FROM SPACEXTBL
           GROUP BY MISSION_OUTCOME;
```

 * db2://vqw61101:***@b0aebb68-94fa-46ec-a1fc-1c999edb6187.c3n41cmd0nqnrk39u98g.databases.appdomain.cloud:31249/bludb
Done.

Out[83]:

| mission_outcome | total_number |
|---|---|
| Failure (in flight) | 1 |
| Success | 99 |
| Success (payload status unclear) | 1 |

# Boosters Carried Maximum Payload

- Names of the booster which have carried the maximum payload mass

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

```
In [84]: %%sql
SELECT DISTINCT BOOSTER_VERSION
FROM SPACEXTBL
WHERE PAYLOAD_MASS__KG_ = (
    SELECT MAX(PAYLOAD_MASS__KG_)
    FROM SPACEXTBL);
```

* db2://vqw61101:***@b0aebb68-94fa-46ec-a1fc-1c999edb6187.c3n41cmd0nqnrk39u98g.databases.appdomain.cloud:31249/bludb
Done.

Out[84]:

| booster_version |
| --- |
| F9 B5 B1048.4 |
| F9 B5 B1048.5 |
| F9 B5 B1049.4 |
| F9 B5 B1049.5 |
| F9 B5 B1049.7 |
| F9 B5 B1051.3 |
| F9 B5 B1051.4 |
| F9 B5 B1051.6 |
| F9 B5 B1056.4 |
| F9 B5 B1058.3 |
| F9 B5 B1060.2 |
| F9 B5 B1060.3 |

# 2015 Launch Records

- Failed landing_outcomes in drone ship, their booster versions, and launch site names for in year 2015

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

```
In [85]:  %%sql
          SELECT Landing__Outcome, BOOSTER_VERSION, LAUNCH_SITE
          FROM SPACEXTBL
          WHERE Landing__Outcome = 'Failure (drone ship)'
              AND YEAR(DATE) = 2015;
```

 * db2://vqw61101:***@b0aebb68-94fa-46ec-a1fc-1c999edb6187.c3n41cmd0nqnrk39u98g.databases.appdomain.cloud:31249/bludb
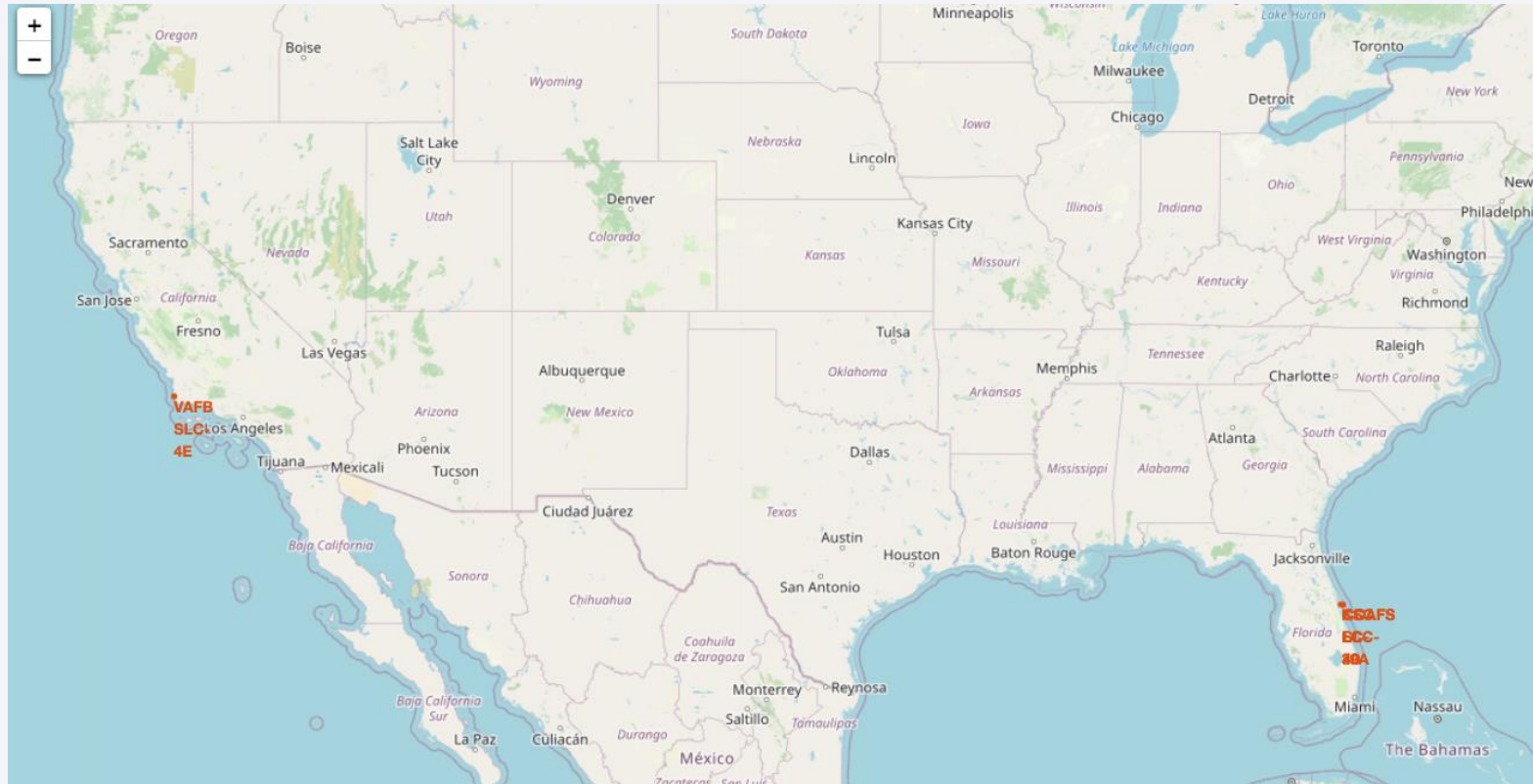Done.

Out[85]:

| landing__outcome | booster_version | launch_site |
|---|---|---|
| Failure (drone ship) | F9 v1.1 B1012 | CCAFS LC-40 |
| Failure (drone ship) | F9 v1.1 B1015 | CCAFS LC-40 |

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

- Ranking 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

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 [86]: %%sql
         SELECT LANDING__OUTCOME, COUNT(LANDING__OUTCOME) AS TOTAL_NUMBER
         FROM SPACEXTBL
         WHERE DATE BETWEEN '2010-06-04' AND '2017-03-20'
         GROUP BY LANDING__OUTCOME
         ORDER BY TOTAL_NUMBER DESC
```

 * db2://vqw61101:***@b0aebb68-94fa-46ec-a1fc-1c999edb6187.c3n41cmd0nqnrk39u98g.databases.appdomain.cloud:31249/bludb
Done.

Out[86]:

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

Section 3

# Launch Sites Proximities Analysis
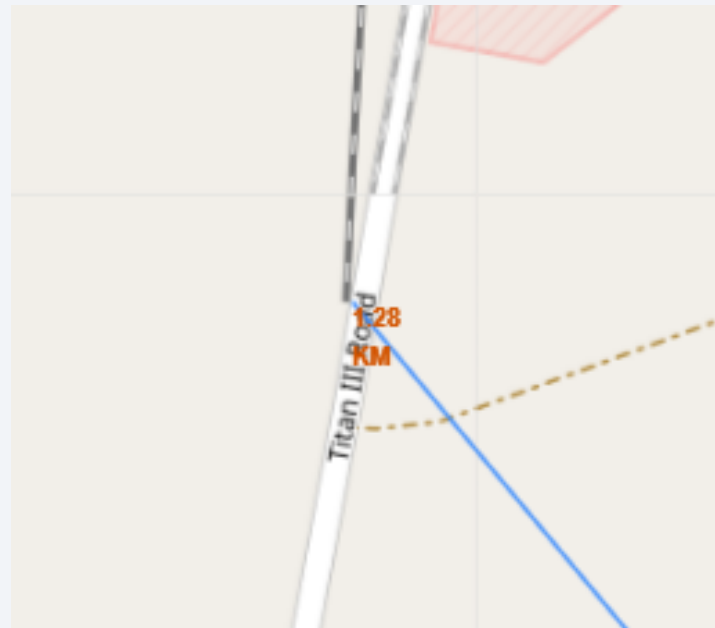
# Location of all the Launch Sites



We can see that all the SpaceX launch sites are located inside the United States

# Markers showing launch sites with color labels

Green markers shows the successful launch and Red markers shows the unsuccessful launch

# Launch Sites Distance to Landmarks



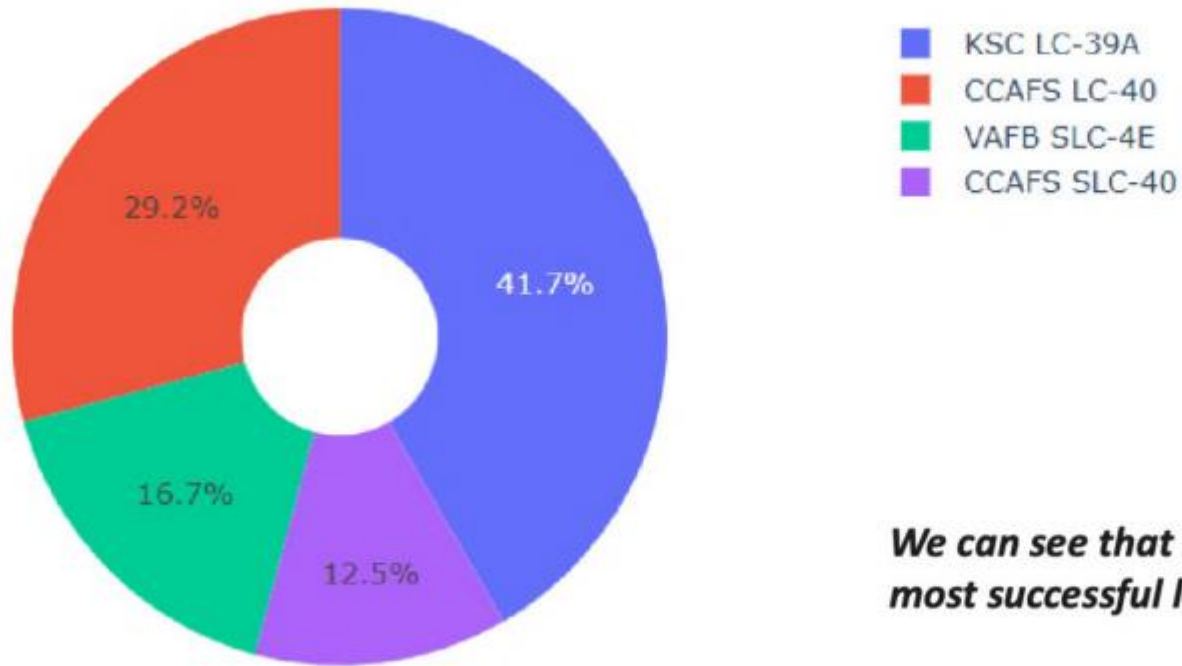Distance lines to the proximities to answer the following questions easily:
- Are launch sites in close proximity to railways? NO
- Are launch sites in close proximity to highways? NO
- Are launch sites in close proximity to coastline? YES
- Do launch sites keep certain distance away from cities? YES

Section 4

# Build a Dashboard
# with Plotly Dash

# The success percentage by each sites



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

29.2%
41.7%
16.7%
12.5%

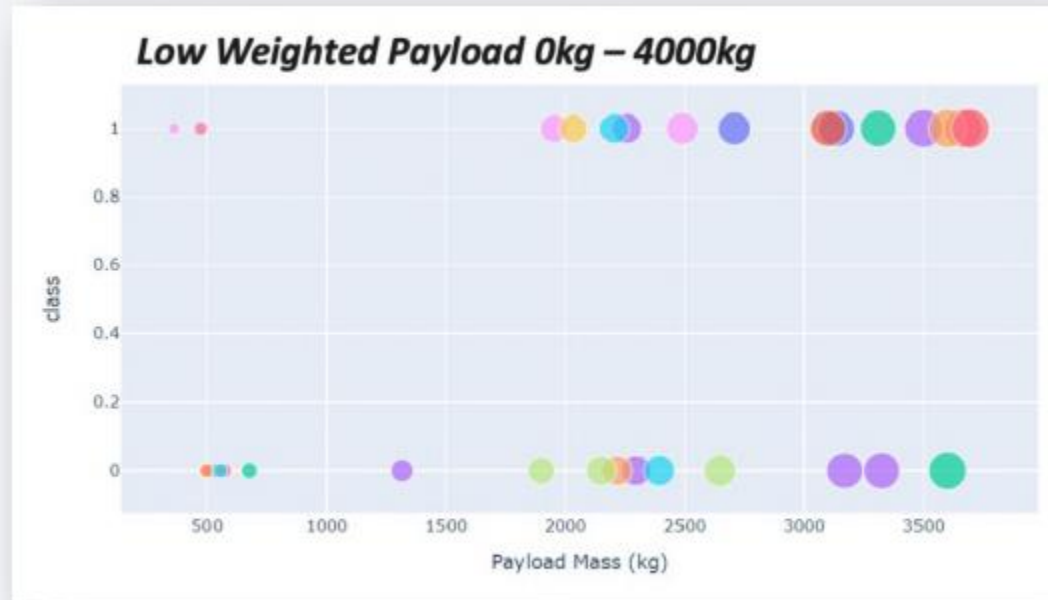*We can see that KSC LC-39A had the most successful launches from all the sites*

# The highest launch-success ratio: KSC LC-39A



23.1%

76.9%

1
0

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

# Payload vs Launch Outcome Scatter Plot

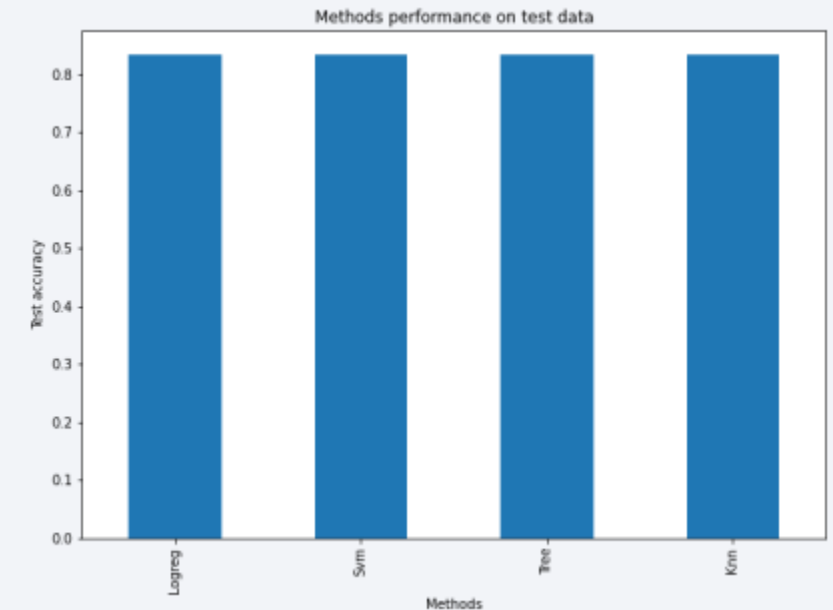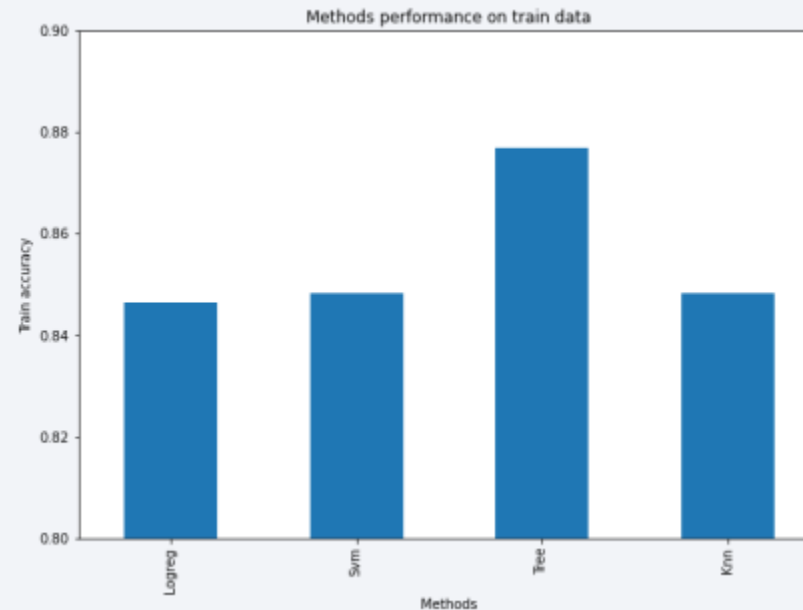We can see that all the success rate for low weighted payload is higher than heavy weighted payload

Section 5

# Predictive Analysis (Classification)

# Classification Accuracy

For accuracy test, all methods performed similar. We could get more test data to decide between them. But if we really need to choose one right now, we will take the decision tree.

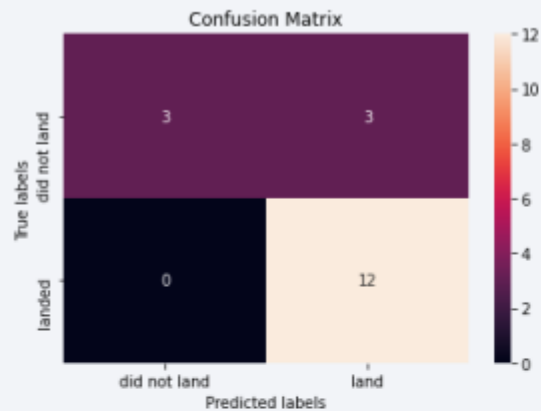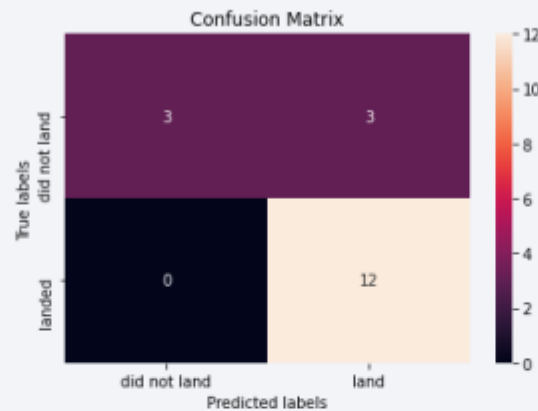|        | Accuracy Train | Accuracy Test |
|--------|----------------|---------------|
| Tree   | 0.876786       | 0.833333      |
| Knn    | 0.848214       | 0.833333      |
| Svm    | 0.848214       | 0.833333      |
| Logreg | 0.846429       | 0.833333      |



**Decision tree best parameters**

```
tuned hyperparameters :(best parameters)  {'criterion': 'entropy', 'max_depth': 12, 'max_features': 'sqrt', 'min_samples_leaf':
4, 'min_samples_split': 2, 'splitter': 'random'}
```
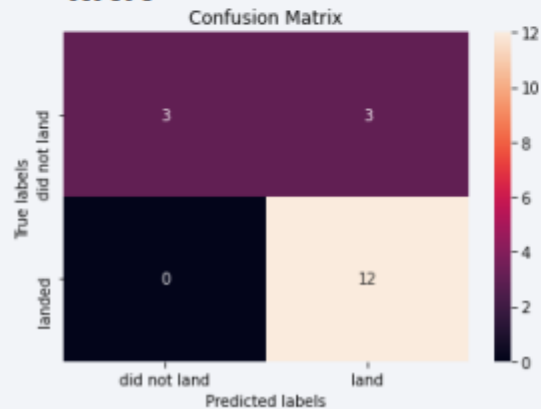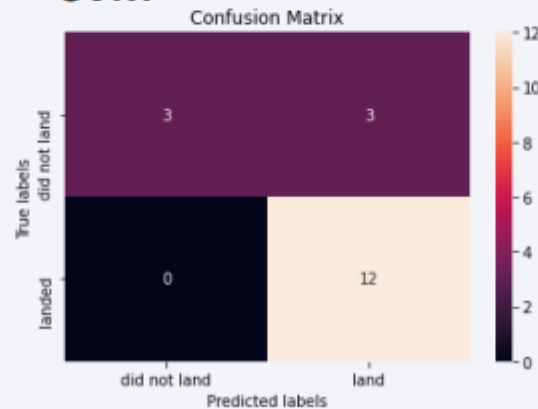
# Confusion Matrix

**Logistic regression**



**Decision Tree**



**kNN**



**SVM**



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

# Conclusions

The success of a mission can be explained by several factors such as the launch site, the orbit and especially the number of previous launches. Indeed, we can assume that there has been a gain in knowledge between launches that allowed to go from a launch failure to a success.

The orbits with the best success rates are GEO, HEO, SSO, ES-L1.

Depending on the orbits, the payload mass can be a criterion to consider for the success of a mission. Some orbits require a light or heavy payload mass. But generally low weighted payloads perform better than the heavy weighted payloads.

With the current data, we cannot explain why some launch sites are better than others (KSC LC-39A is the best launch site). To get an answer to this problem, we could obtain atmospheric or other relevant data.

For this dataset, we choose the Decision Tree Algorithm as the best model even if the test accuracy between all the models used is identical. We choose Decision Tree Algorithm because it has a better train accuracy.

Thank you!