



# Exploring Space with Data Science

Manish Singh

22/11/2022

# OUTLINE

---



- Executive Summary
- Introduction
- Methodology
- Results
  - Visualization – Charts
  - Dashboard
- Discussion
  - Findings & Implications
- 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 Data Visualization
  - Interactive Visual Analytics with Folium
  - Machine Learning Prediction
- Summary of all results
  - Exploratory Data Analysis result
  - Interactive analytics in screenshots
  - Predictive Analytics result from Machine Learning Lab

# INTRODUCTION

---



- SpaceX is a revolutionary company who has disrupted the space industry by offering rocket launches specifically Falcon 9 as low as 62 million dollars; while other providers cost upward of 165 million dollars each. Most of this saving thanks to SpaceX's astounding idea to reuse the first stage of the launch by re-land the rocket to be used on the next mission. Repeating this process will make the price down even further. As a data scientist of a startup rivaling SpaceX, the goal of this project is to create the machine learning pipeline to predict the landing outcome of the first stage in the future. This project is crucial in identifying the right price to bid against SpaceX for a rocket launch.
- The problems included:
  - Identifying all factors that influence the landing outcome.
  - The relationship between each variable and how it is affecting the outcome.
  - The best condition needed to increase the probability of successful landing.

# METHODOLOGY

---



## Executive Summary

- Data collection methodology:
  - Data was collected using SpaceX REST API and web scrapping from Wikipedia
- Perform data wrangling
  - Data was processed 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

---

Data collection is the process of gathering and measuring information on targeted variables in an established system, which then enables one to answer relevant questions and evaluate outcomes. As mentioned, the dataset was collected by REST API and Web Scrapping from Wikipedia

For REST API, its started by using the get request. Then, we decoded the response content as Json and turn it into a pandas dataframe using `json_normalize()`. We then cleaned the data, checked for missing values and fill with whatever needed.

For web scrapping, we will use the BeautifulSoup to extract the launch records as HTML table, parse the table and convert it to a pandas dataframe for further analysis

# DATA COLLECTION – SpaceX API

Get request for rocket launch data using AP

Use json\_normalize method to convert json result to dataframe

Performed data cleaning and filling the missing value

```
spacex_url="https://api.spacexdata.com/v4/launches/past"
```

```
response = requests.get(spacex_url)
```

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

```
# Lets take a subset of our dataframe keeping only the features we want a  
nd the flight number, and date_utc.  
data = data[['rocket', 'payloads', 'launchpad', 'cores', 'flight_number',  
'date_utc']]
```

```
# We will remove rows with multiple cores because those are falcon rocket  
s with 2 extra rocket boosters and rows that have multiple payloads in a  
single rocket.  
data = data[data['cores'].map(len)==1]  
data = data[data['payloads'].map(len)==1]
```

```
# Since payloads and cores are lists of size 1 we will also extract the s  
ingle value in the list and replace the feature.  
data['cores'] = data['cores'].map(lambda x : x[0])  
data['payloads'] = data['payloads'].map(lambda x : x[0])
```

```
# We also want to convert the date_utc to a datetime datatype and then ex  
tracting the date leaving the time  
data['date'] = pd.to_datetime(data['date_utc']).dt.date
```

```
# Using the date we will restrict the dates of the launches  
data = data[data['date'] <= datetime.date(2020, 11, 13)]
```



# DATA COLLECTION – Scraping

Request the Falcon9 Launch Wiki page from url

Create a BeautifulSoup from the HTML response

Extract all column/variable names from the HTML header

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

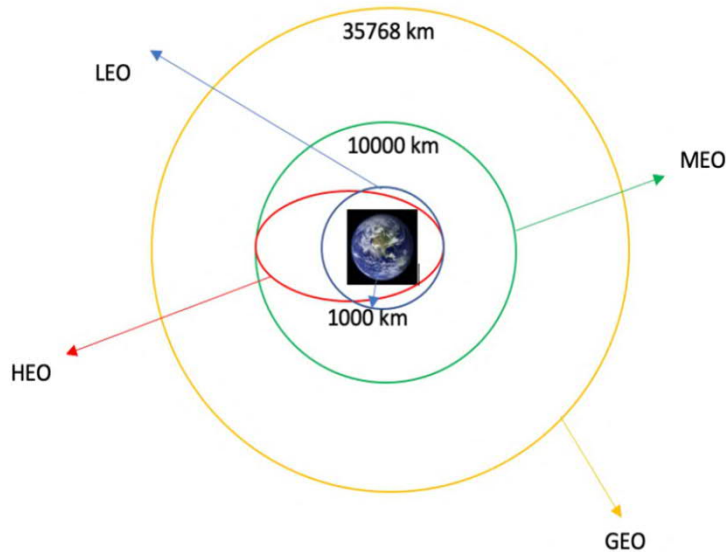
```
# Use BeautifulSoup() to create a BeautifulSoup object from a response text content
soup = BeautifulSoup(data, 'html.parser')
```

```
extracted_row = 0
#Extract each table
for table_number, table in enumerate(soup.find_all('table', "wikitable plainrowheaders collapsible")):
    # get table row
    for rows in table.find_all("tr"):
        #check to see if first table heading is as number corresponding to launch a number
        if rows.th:
            if rows.th.string:
                flight_number=rows.th.string.strip()
                flag=flight_number.isdigit()
            else:
                flag=False
```



# DATA WRANGLING

---

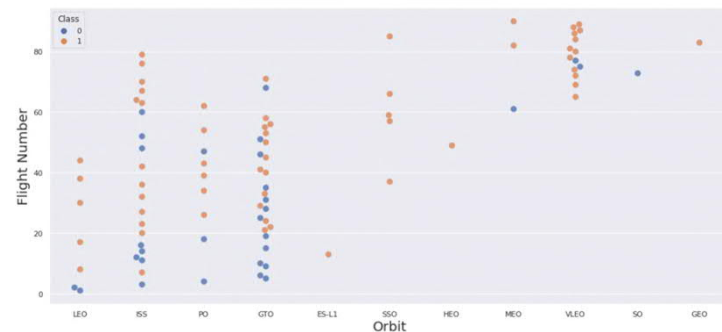
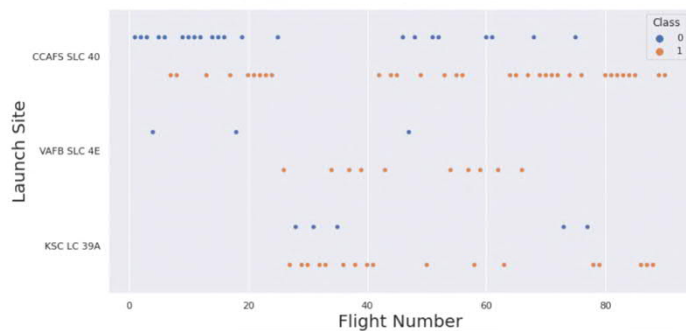


Data Wrangling is the process of cleaning and unifying messy and complex data sets for easy access and Exploratory Data Analysis (EDA).

We will first calculate the number of launches on each site, then calculate the number and occurrence of mission outcome per orbit type.

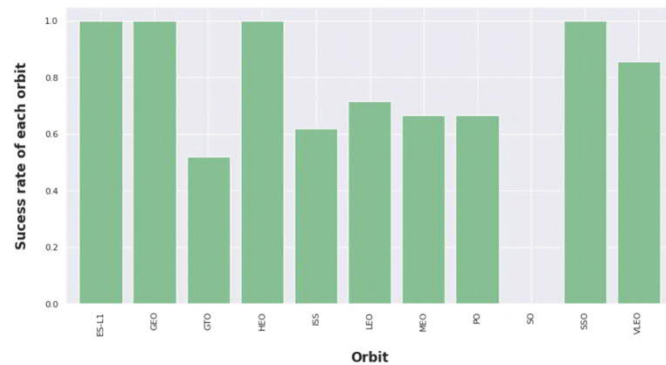
We then create a landing outcome label from the outcome column. This will make it easier for further analysis, visualization, and ML. Lastly, we will export the result to a CSV.

# EDA WITH DATA VISUALIZATION



- We first started by using scatter graph to find the relationship between the attributes such as between:
  - Payload and Flight Number.
  - Flight Number and Launch Site.
  - Payload and Launch Site.
  - Flight Number and Orbit Type.
  - Payload and Orbit Type.
- Scatter plots show dependency of attributes on each other. Once a pattern is determined from the graphs. It's very easy to see which factors affecting the most to the success of the landing outcomes.

# EDA WITH DATA VISUALIZATION

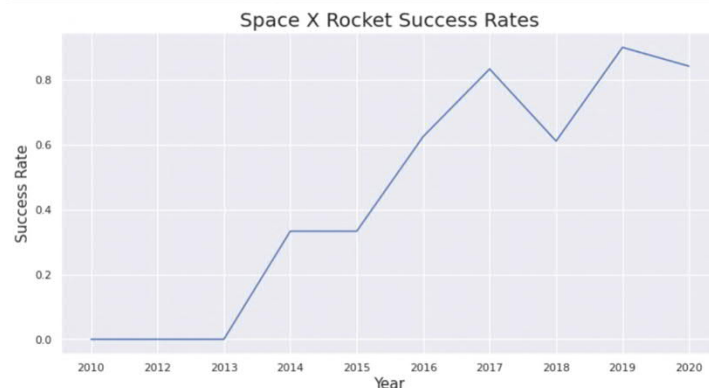


Once we get a hint of the relationships using scatter plot. We will then use further visualization tools such as bar graph and line plots graph for further analysis.

Bar graphs is one of the easiest way to interpret the relationship between the attributes. In this case, we will use the bar graph to determine which orbits have the highest probability of success.

We then use the line graph to show a trends or pattern of the attribute over time which in this case, is used for see the launch success yearly trend.

We then use Feature Engineering to be used in success prediction in the future module by created the dummy variables to categorical columns.



# EDA WITH SQL

---

Using SQL, we had performed many queries to get better understanding of the dataset, Ex:

- Displaying the names of the launch sites.
- Displaying 5 records where launch sites begin with the string 'CCA'.
- Displaying the total payload mass carried by booster launched by NASA (CRS).
- Displaying the average payload mass carried by booster version F9 v1.1.
- Listing the date when the first successful landing outcome in ground pad was achieved.
- Listing the names of the boosters which have success in drone ship and have payload mass greater than 4000 but less than 6000.
- Listing the total number of successful and failure mission outcomes.
- Listing the names of the booster\_versions which have carried the maximum payload mass.
- Listing the failed landing\_outcomes in drone ship, their booster versions, and launch sites names for in year 2015.
- Rank the count of landing outcomes or success between the date 2010-06-04 and 2017-03-20, in descending order

# EDA WITH SQL

LOAD DATA

Source Target Define Finalize

You are loading the file **Spacex.csv**

Select a load target

Schema	Table
<input type="text" value="Find a schema"/>	<input type="text" value="Find a table in QWP24235"/>
AUDIT	ANNUAL_CROP_DATA
DRZINST1	BOARD
ERRORSCHEMA	BOOKSHOP

LOAD DATA

Source Target Define Finalize

You are loading the file **Spacex.csv** into **QWP24235.SPACEXTBL**

Code page (character encoding): 1208 (UTF-8) Separator: Header in first row: ☒ Time & date format: Detect data types: ☐

Date format: DD-MM-YYYY Time format: HH:MM:SS Timestamp format: DD-MM-YYYY HH:MM:SS

LAUNCH_SITE	PAYLOAD	PAYLOAD_MASS_KG	ORBIT	CUSTOMER
CCAFS LC-10	Dragon Spacecraft Qualification Unit			SpaceX

SQL

Source Target Define Finalize

You are loading the file **Spacex (2).csv** into **SRW76180.SPACEXTBL**

Code page (character encoding): 1208 (UTF-8) Separator: Header in first row: ☒ Time & date format:

Date format: DD-MM-YYYY Time format: HH:MM:SS Timestamp format: DD-MM-YYYY HH:MM:SS

DATE	TIME_UTC	BOOSTER_VERSION	LAUNCH_SITE	PAYLOAD	PAYLOAD_SMALL
04-06-2010	18:45:00	F9 v1.0 B0003	CCAFS LC-40	Dragon Spacecraft Qualification Unit	0
08-12-2010	15:43:00	F9 v1.0 B0004	CCAFS LC-40	Dragon demo flight C1, two CubeSats, barrel of Brounche cheese	0

# BUILD AN INTERACTIVE MAP WITH FOLIUM

---

To visualize the launch data into an interactive map. We took the latitude and longitude coordinates at each launch site and added a circle marker around each launch site with a label of the name of the launch site.

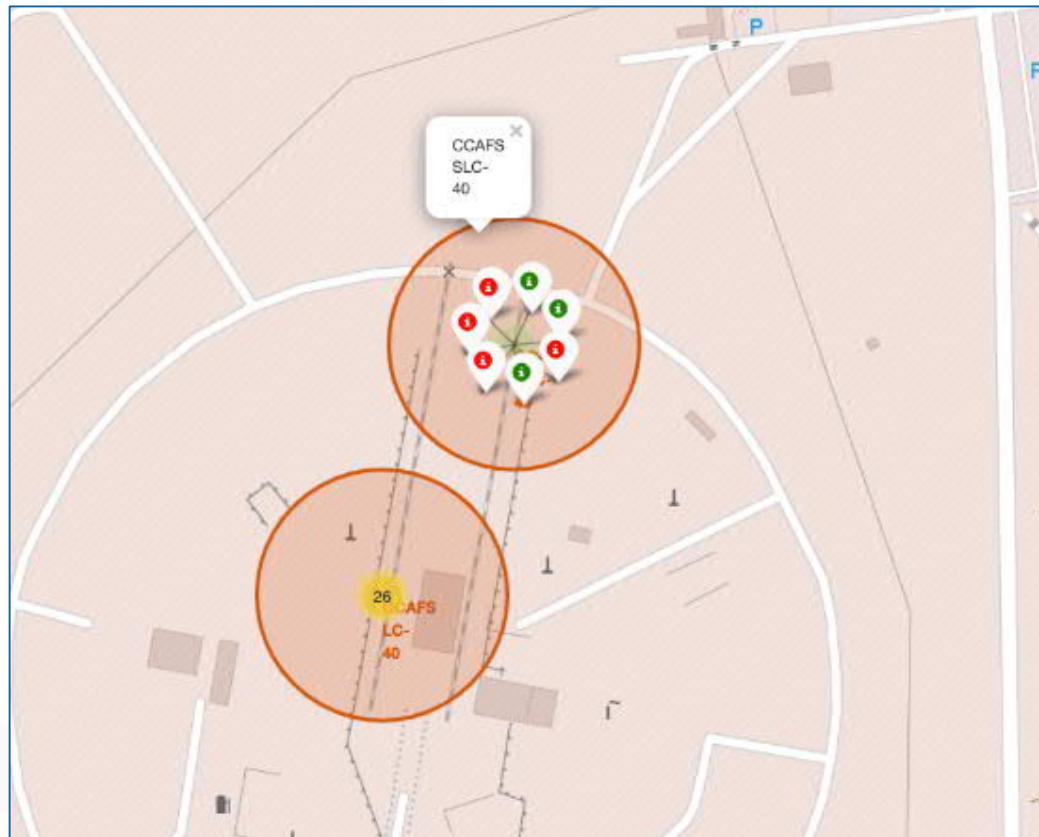
We then assigned the dataframe `launch_outcomes(failure,success)` to classes 0 and 1 with **Red** and **Green** markers on the map in `MarkerCluster()`.

We then used the Haversine's formula to calculate the distance of the launch sites to various landmarks to find answers to the questions of:

- How close the launch sites with railways, highways and coastlines?
- How close the launch sites with nearby cities?

# BUILD AN INTERACTIVE MAP WITH FOLIUM

---





# BUILD A DASHBOARD WITH PLOTLY DASH

- We built an interactive dashboard with Plotly dash which allowing the user to play around with the data as they need.
- We plotted pie charts showing the total launches by a certain sites.
- We then plotted scatter graph showing the relationship with Outcome and Payload Mass (Kg) for the different booster version.

```
1 import pandas as pd
2 import dash
3 import dash_html_components as html
4 import dash_core_components as dcc
5 from dash.dependencies import Input, Output, Callback
6 import plotly.express as px
7
8 # Load data
9 df = pd.read_csv('https://raw.githubusercontent.com/plotly/datasets/master/SpacexLaunches.csv')
10
11 # Filter data for the first 100 launches
12 df = df[:100]
13
14 # Create a Dash app
15 app = dash.Dash(__name__)
16 server = app.server
17
18 # Define the layout
19 app.layout = html.Div([
20     dcc.Graph(id='total-launches'),
21     dcc.Graph(id='outcome-payload'),
22     dcc.Graph(id='payload-mass'),
23 ])
24
25 # Define the callbacks
26 @app.callback([
27     Output('total-launches', 'figure'),
28     Output('outcome-payload', 'figure'),
29     Output('payload-mass', 'figure')
30 ], [
31     Input('total-launches', 'figure'),
32     Input('outcome-payload', 'figure'),
33     Input('payload-mass', 'figure')
34 ])
35 def update_figures(fig1, fig2, fig3):
36     # Update the total launches figure
37     fig1 = px.pie(df, names='site', title='Total Launches by Site')
38
39     # Update the outcome-payload figure
40     fig2 = px.scatter(df, x='payload_mass_kg', y='outcome', title='Outcome vs Payload Mass (Kg)')
41
42     # Update the payload-mass figure
43     fig3 = px.scatter(df, x='payload_mass_kg', y='payload_mass_kg', title='Payload Mass (Kg)')
44
45     return fig1, fig2, fig3
46
47 # Run the app
48 if __name__ == '__main__':
49     app.run_server(debug=True)
```

# 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.

# Results

---

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

- Exploratory data analysis results
- Interactive analytics demo in screenshots
- Predictive analysis results