# SpaceX Falcon 9 First Stage Landing Prediction

## Assignment: Exploring and Preparing Data

Estimated time needed: **70** minutes

In this assignment, we will predict if the Falcon 9 first stage will land successfully. 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 due to the fact that SpaceX can reuse the first stage.

In this lab, you will perform Exploratory Data Analysis and Feature Engineering.

Falcon 9 first stage will land successfully



Several examples of an unsuccessful landing are shown here:



Most unsuccessful landings are planned. Space X performs a controlled landing in the oceans.

## Objectives

Perform exploratory Data Analysis and Feature Engineering using `Pandas` and `Matplotlib`

- Exploratory Data Analysis
- Preparing Data Feature Engineering

### Import Libraries and Define Auxiliary Functions

We will import the following libraries the lab

```
In [ ]:  import piplite
         await piplite.install(['numpy'])
         await piplite.install(['pandas'])
         await piplite.install(['seaborn'])
```

```
In [ ]:  # pandas is a software library written for the Python programming language for data manipulation and analysis.
         import pandas as pd
         #NumPy is a library for the Python programming language, adding support for large, multi-dimensional arrays and matrices, alon
         import numpy as np
         # Matplotlib is a plotting library for python and pyplot gives us a MatLab like plotting framework. We will use this in our pl
```

```
import matplotlib.pyplot as plt
#Seaborn is a Python data visualization library based on matplotlib. It provides a high-level interface for drawing attractive
import seaborn as sns
```

# Exploratory Data Analysis

First, let's read the SpaceX dataset into a Pandas dataframe and print its summary

```
from js import fetch
import io

URL = "https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBM-DS0321EN-SkillsNetwork/datasets/dataset_part_2.c
resp = await fetch(URL)
dataset_part_2_csv = io.BytesIO((await resp.arrayBuffer()).to_py())
df=pd.read_csv(dataset_part_2_csv)
df.head(5)
```
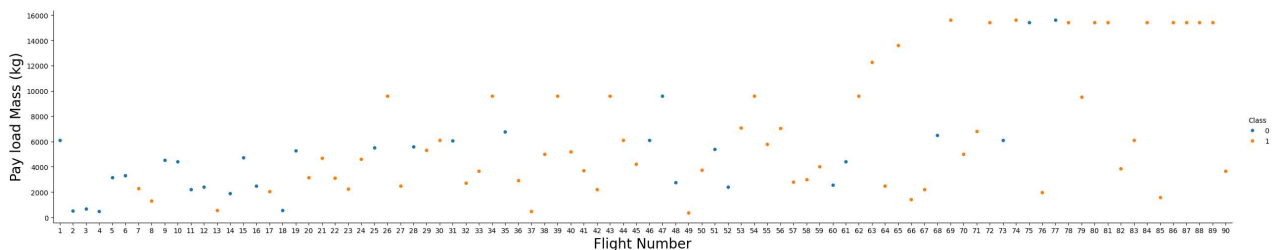
Out[ ]:

| | FlightNumber | Date | BoosterVersion | PayloadMass | Orbit | LaunchSite | Outcome | Flights | GridFins | Reused | Legs | LandingPad | Block | ReusedCou |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 2010-06-04 | Falcon 9 | 6104.959412 | LEO | CCAFS SLC 40 | None None | 1 | False | False | False | NaN | 1.0 | |
| 1 | 2 | 2012-05-22 | Falcon 9 | 525.000000 | LEO | CCAFS SLC 40 | None None | 1 | False | False | False | NaN | 1.0 | |
| 2 | 3 | 2013-03-01 | Falcon 9 | 677.000000 | ISS | CCAFS SLC 40 | None None | 1 | False | False | False | NaN | 1.0 | |
| 3 | 4 | 2013-09-29 | Falcon 9 | 500.000000 | PO | VAFB SLC 4E | False Ocean | 1 | False | False | False | NaN | 1.0 | |
| 4 | 5 | 2013-12-03 | Falcon 9 | 3170.000000 | GTO | CCAFS SLC 40 | None None | 1 | False | False | False | NaN | 1.0 | |

First, let's try to see how the `FlightNumber` (indicating the continuous launch attempts.) and `Payload` variables would affect the launch outcome.

We can plot out the `FlightNumber` vs. `PayloadMass` and overlay the outcome of the launch. We see that as the flight number increases, the first stage is more likely to land successfully. The payload mass also appears to be a factor; even with more massive payloads, the first stage often returns successfully.

```
sns.catplot(y="PayloadMass", x="FlightNumber", hue="Class", data=df, aspect = 5)
plt.xlabel("Flight Number",fontsize=20)
plt.ylabel("Pay load Mass (kg)",fontsize=20)
plt.show()
```
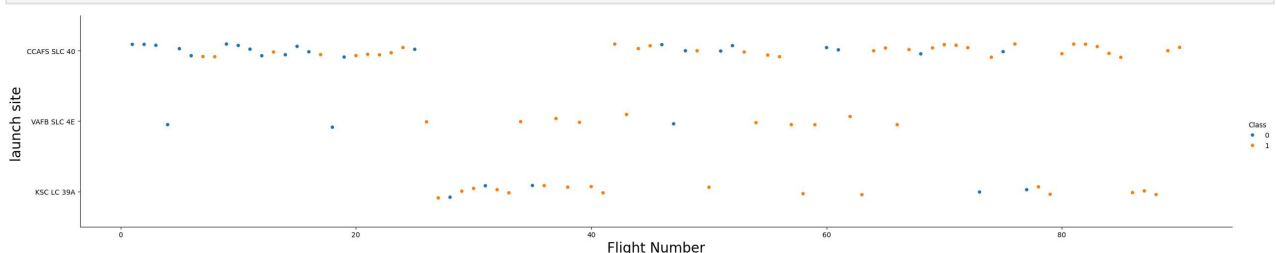


Next, let's drill down to each site visualize its detailed launch records.

## TASK 1: Visualize the relationship between Flight Number and Launch Site

Use the function `catplot` to plot `FlightNumber` vs `LaunchSite`, set the parameter `x` parameter to `FlightNumber`,set the `y` to `Launch Site` and set the parameter `hue` to `'class'`

```
# Plot a scatter point chart with x axis to be Flight Number and y axis to be the launch site, and hue to be the class value

sns.catplot(y="LaunchSite", x="FlightNumber", hue="Class", data=df, aspect = 5)
plt.xlabel("Flight Number",fontsize=20)
plt.ylabel("launch site ",fontsize=20)
plt.show()
```
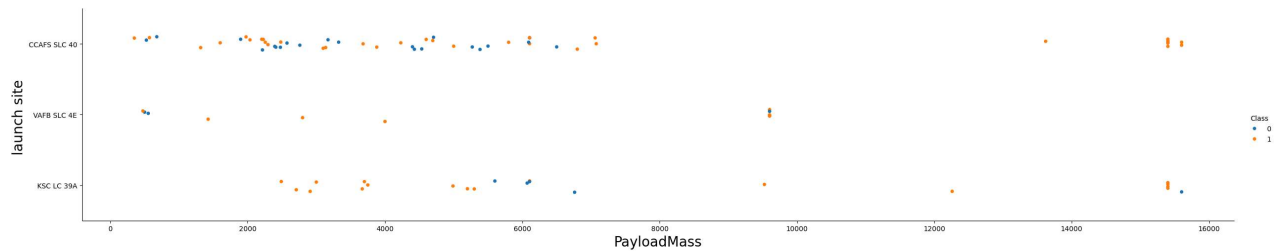


Now try to explain the patterns you found in the Flight Number vs. Launch Site scatter point plots.

## TASK 2: Visualize the relationship between Payload Mass and Launch Site

We also want to observe if there is any relationship between launch sites and their payload mass.

```
In [9]:  # Plot a scatter point chart with x axis to be Pay Load Mass (kg) and y axis to be the launch site, and hue to be the class va
         sns.catplot(y="LaunchSite", x="PayloadMass", hue="Class", data=df, aspect = 5)
         plt.xlabel("PayloadMass",fontsize=20)
         plt.ylabel("launch site ",fontsize=20)
         plt.show()
```



Now if you observe Payload Mass Vs. Launch Site scatter point chart you will find for the VAFB-SLC launchsite there are no rockets launched for heavypayload mass(greater than 10000).

## TASK 3: Visualize the relationship between success rate of each orbit type

Next, we want to visually check if there are any relationship between success rate and orbit type.

Let's create a `bar chart` for the sucess rate of each orbit

```
In [10]:  # HINT use groupby method on Orbit column and get the mean of Class column

          # Calculate the success rate for each orbit
          success_rate = df.groupby('Orbit')['Class'].mean()

          # Sort the values from highest to lowest success rate
          success_rate = success_rate.sort_values(ascending=False)

          # Create a bar chart
          plt.figure(figsize=(10,6))
          success_rate.plot(kind='bar', color='skyblue', edgecolor='black')

          # Add titles and labels
          plt.title('Success Rate by Orbit', fontsize=14)
          plt.xlabel('Orbit', fontsize=12)
          plt.ylabel('Success Rate', fontsize=12)
          plt.xticks(rotation=45)
          plt.grid(axis='y', linestyle='--', alpha=0.7)

          # Show the chart
          plt.show()
```
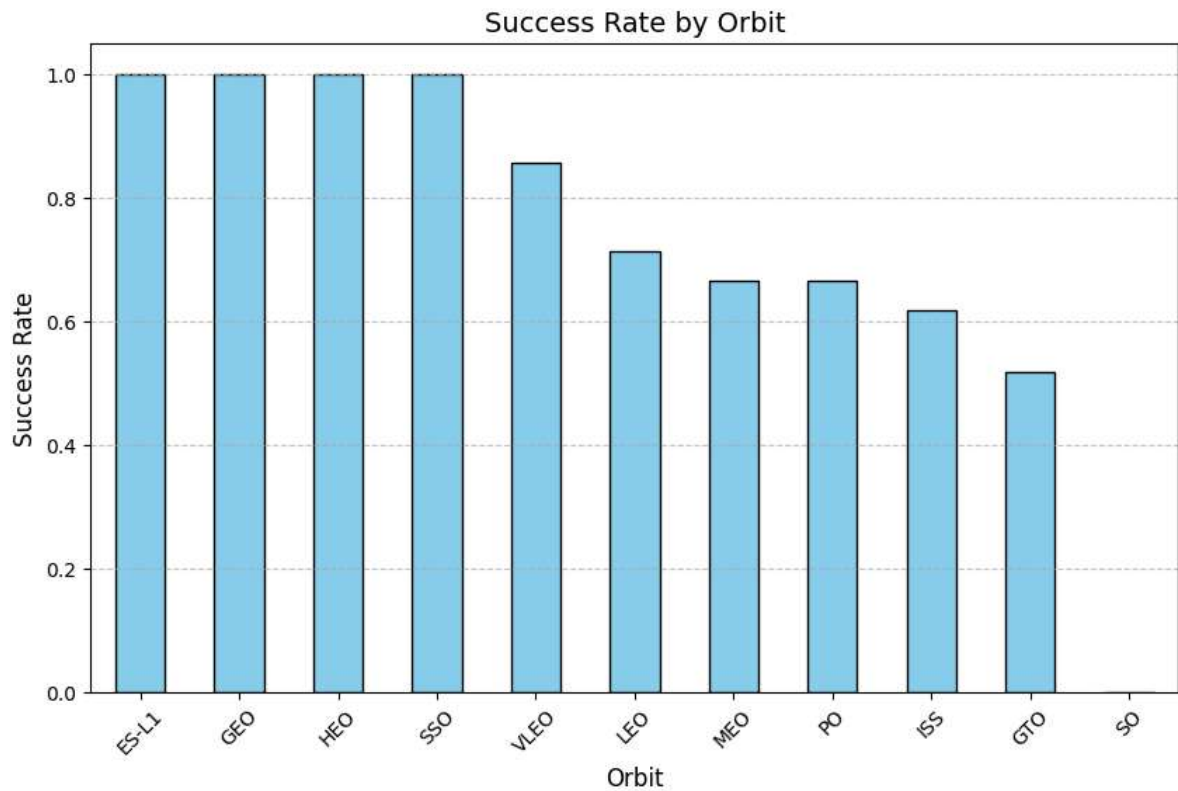
## Success Rate by Orbit



Analyze the plotted bar chart to identify which orbits have the highest success rates.

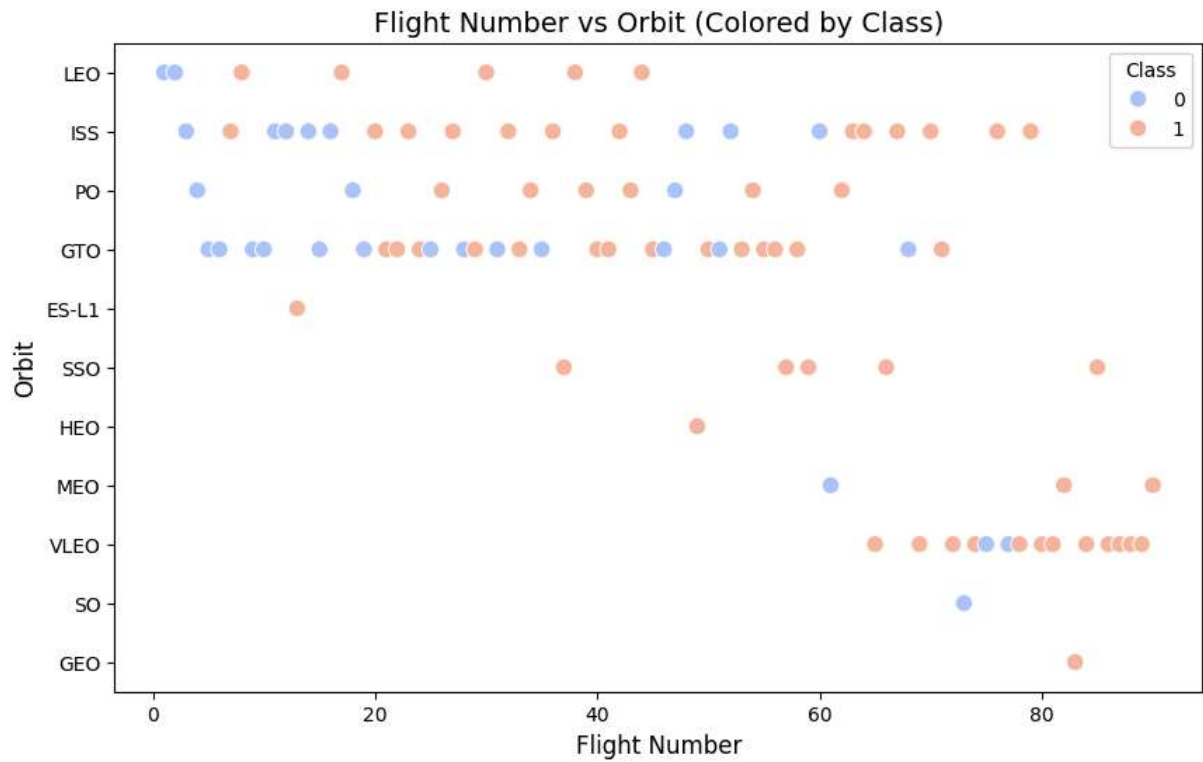## TASK 4: Visualize the relationship between FlightNumber and Orbit type

For each orbit, we want to see if there is any relationship between FlightNumber and Orbit type.

In [11]:
```python
# Plot a scatter point chart with x axis to be FlightNumber and y axis to be the Orbit, and hue to be the class value

# Create a scatter plot
plt.figure(figsize=(10,6))
sns.scatterplot(data=df, x='FlightNumber', y='Orbit', hue='Class', palette='coolwarm', s=80)

# Add titles and labels
plt.title('Flight Number vs Orbit (Colored by Class)', fontsize=14)
plt.xlabel('Flight Number', fontsize=12)
plt.ylabel('Orbit', fontsize=12)
plt.legend(title='Class', loc='best')

# Show the plot
plt.show()
```

Flight Number vs Orbit (Colored by Class)

You can observe that in the LEO orbit, success seems to be related to the number of flights. Conversely, in the GTO orbit, there appears to be no relationship between flight number and success.
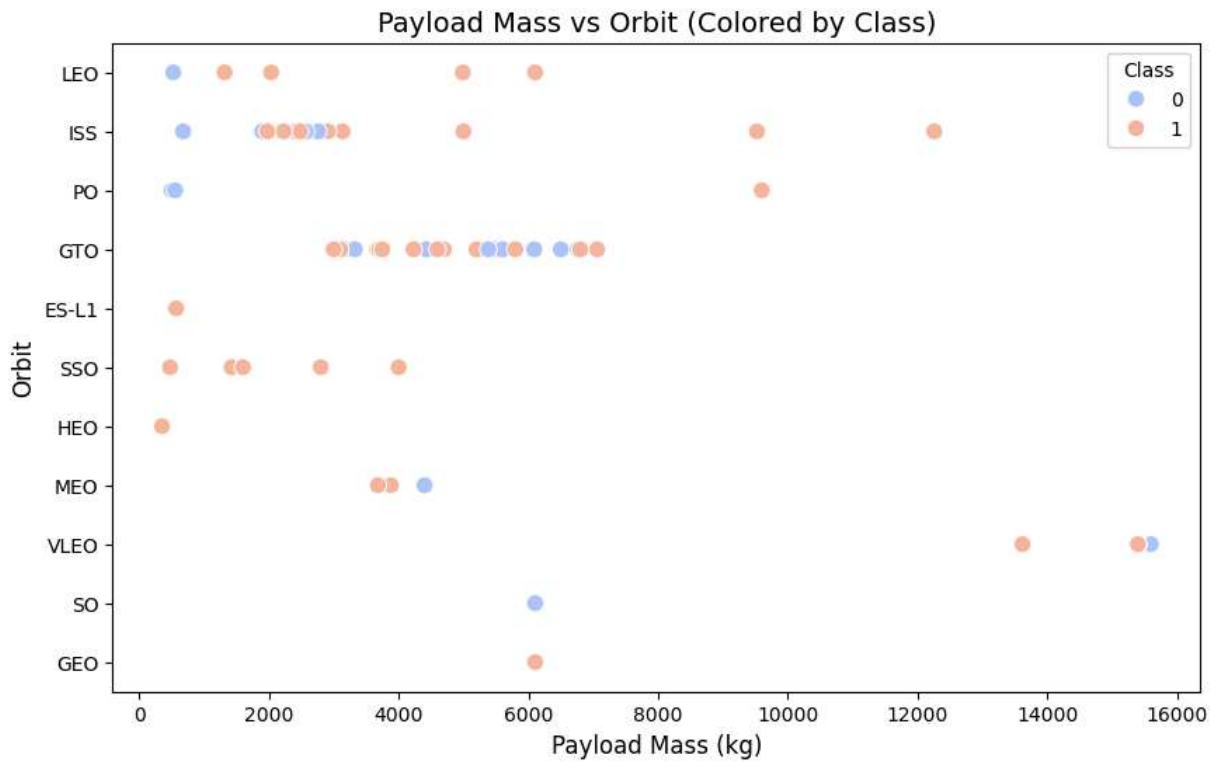
## TASK 5: Visualize the relationship between Payload Mass and Orbit type

Similarly, we can plot the Payload Mass vs. Orbit scatter point charts to reveal the relationship between Payload Mass and Orbit type

```
In [12]:  # Plot a scatter point chart with x axis to be Payload Mass and y axis to be the Orbit, and hue to be the class value
          # Create a scatter plot of Payload Mass vs Orbit, colored by Class
          plt.figure(figsize=(10,6))
          sns.scatterplot(data=df, x='PayloadMass', y='Orbit', hue='Class', palette='coolwarm', s=80)

          # Add titles and labels
          plt.title('Payload Mass vs Orbit (Colored by Class)', fontsize=14)
          plt.xlabel('Payload Mass (kg)', fontsize=12)
          plt.ylabel('Orbit', fontsize=12)
          plt.legend(title='Class', loc='best')

          # Show the plot
          plt.show()
```

Payload Mass vs Orbit (Colored by Class)

With heavy payloads the successful landing or positive landing rate are more for Polar,LEO and ISS.

However, for GTO, it's difficult to distinguish between successful and unsuccessful landings as both outcomes are present.

## TASK 6: Visualize the launch success yearly trend

You can plot a line chart with x axis to be `Year` and y axis to be average success rate, to get the average launch success trend.

The function will help you get the year from the date:

```
In [13]:  # A function to Extract years from the date
          year=[]
          def Extract_year():
              for i in df["Date"]:
                  year.append(i.split("-")[0])
              return year
          Extract_year()
          df['Date'] = year
          df.head()
```

Out[13]:

| | FlightNumber | Date | BoosterVersion | PayloadMass | Orbit | LaunchSite | Outcome | Flights | GridFins | Reused | Legs | LandingPad | Block | ReusedCoun |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 2010 | Falcon 9 | 6104.959412 | LEO | CCAFS SLC 40 | None None | 1 | False | False | False | NaN | 1.0 | 0 |
| 1 | 2 | 2012 | Falcon 9 | 525.000000 | LEO | CCAFS SLC 40 | None None | 1 | False | False | False | NaN | 1.0 | 0 |
| 2 | 3 | 2013 | Falcon 9 | 677.000000 | ISS | CCAFS SLC 40 | None None | 1 | False | False | False | NaN | 1.0 | 0 |
| 3 | 4 | 2013 | Falcon 9 | 500.000000 | PO | VAFB SLC 4E | False Ocean | 1 | False | False | False | NaN | 1.0 | 0 |
| 4 | 5 | 2013 | Falcon 9 | 3170.000000 | GTO | CCAFS SLC 40 | None None | 1 | False | False | False | NaN | 1.0 | 0 |

```
In [14]:  # Plot a line chart with x axis to be the extracted year and y axis to be the success rate

          # Make sure the Date column is in datetime format
          df['Date'] = pd.to_datetime(df['Date'], errors='coerce')

          # Extract the year from the Date column
          df['Year'] = df['Date'].dt.year

          # Calculate the success rate per year
          yearly_success_rate = df.groupby('Year')['Class'].mean()

          # Create a line chart
          plt.figure(figsize=(10,6))
          plt.plot(yearly_success_rate.index, yearly_success_rate.values, marker='o', linestyle='-', color='royalblue')
```
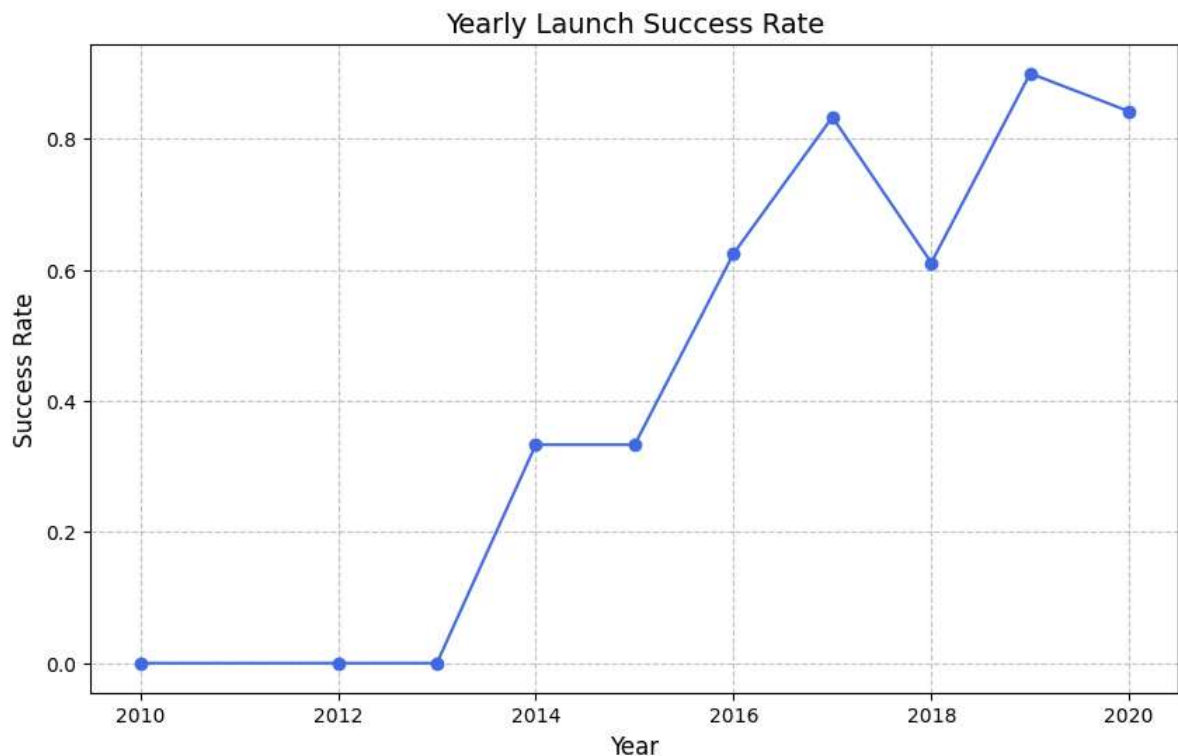
```
# Add titles and labels
plt.title('Yearly Launch Success Rate', fontsize=14)
plt.xlabel('Year', fontsize=12)
plt.ylabel('Success Rate', fontsize=12)
plt.grid(True, linestyle='--', alpha=0.7)

# Show the chart
plt.show()
```



you can observe that the sucess rate since 2013 kept increasing till 2020

# Features Engineering

By now, you should obtain some preliminary insights about how each important variable would affect the success rate, we will select the features that will be used in success prediction in the future module.

In [15]: 
```
features = df[['FlightNumber', 'PayloadMass', 'Orbit', 'LaunchSite', 'Flights', 'GridFins', 'Reused', 'Legs', 'LandingPad', 'B
features.head()
```

Out[15]:

| | FlightNumber | PayloadMass | Orbit | LaunchSite | Flights | GridFins | Reused | Legs | LandingPad | Block | ReusedCount | Serial |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 1 | 6104.959412 | LEO | CCAFS SLC 40 | 1 | False | False | False | NaN | 1.0 | 0 | B0003 |
| **1** | 2 | 525.000000 | LEO | CCAFS SLC 40 | 1 | False | False | False | NaN | 1.0 | 0 | B0005 |
| **2** | 3 | 677.000000 | ISS | CCAFS SLC 40 | 1 | False | False | False | NaN | 1.0 | 0 | B0007 |
| **3** | 4 | 500.000000 | PO | VAFB SLC 4E | 1 | False | False | False | NaN | 1.0 | 0 | B1003 |
| **4** | 5 | 3170.000000 | GTO | CCAFS SLC 40 | 1 | False | False | False | NaN | 1.0 | 0 | B1004 |

## TASK 7: Create dummy variables to categorical columns

Use the function `get_dummies` and `features` dataframe to apply OneHotEncoder to the column `Orbits`, `LaunchSite`, `LandingPad`, and `Serial`. Assign the value to the variable `features_one_hot`, display the results using the method head. Your result dataframe must include all features including the encoded ones.

In [16]: 
```
# HINT: Use get_dummies() function on the categorical columns

# Apply one-hot encoding to the categorical columns
features_one_hot = pd.get_dummies(features,
                                  columns=['Orbit', 'LaunchSite', 'LandingPad', 'Serial'])

# Display the first 5 rows of the encoded DataFrame
features_one_hot.head()
```

| | FlightNumber | PayloadMass | Flights | GridFins | Reused | Legs | Block | ReusedCount | Orbit_ES-L1 | Orbit_GEO | ... | Serial_B1048 | Serial_B1049 | Serial_B1( |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 1 | 6104.959412 | 1 | False | False | False | 1.0 | 0 | False | False | ... | False | False | Fa |
| **1** | 2 | 525.000000 | 1 | False | False | False | 1.0 | 0 | False | False | ... | False | False | Fa |
| **2** | 3 | 677.000000 | 1 | False | False | False | 1.0 | 0 | False | False | ... | False | False | Fa |
| **3** | 4 | 500.000000 | 1 | False | False | False | 1.0 | 0 | False | False | ... | False | False | Fa |
| **4** | 5 | 3170.000000 | 1 | False | False | False | 1.0 | 0 | False | False | ... | False | False | Fa |

5 rows × 80 columns

## TASK 8: Cast all numeric columns to `float64`

Now that our `features_one_hot` dataframe only contains numbers, cast the entire dataframe to variable type `float64`

In [17]:
```python
# HINT: use astype function
# Convert the entire DataFrame to float64 type
features_one_hot = features_one_hot.astype('float64')

# Display the first 5 rows to confirm the change
features_one_hot.head()
```

Out[17]:

| | FlightNumber | PayloadMass | Flights | GridFins | Reused | Legs | Block | ReusedCount | Orbit_ES-L1 | Orbit_GEO | ... | Serial_B1048 | Serial_B1049 | Serial_B1( |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 1.0 | 6104.959412 | 1.0 | 0.0 | 0.0 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 | |
| **1** | 2.0 | 525.000000 | 1.0 | 0.0 | 0.0 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 | |
| **2** | 3.0 | 677.000000 | 1.0 | 0.0 | 0.0 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 | |
| **3** | 4.0 | 500.000000 | 1.0 | 0.0 | 0.0 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 | |
| **4** | 5.0 | 3170.000000 | 1.0 | 0.0 | 0.0 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 | |

5 rows × 80 columns

In [19]:
```python
features_one_hot.to_csv('dataset_part_3.csv', index=False)
```

We can now export it to a **CSV** for the next section,but to make the answers consistent, in the next lab we will provide data in a pre-selected date range.

```
features_one_hot.to_csv('dataset_part\_3.csv', index=False)
```

## Authors

Pratiksha Verma