

# SpaceX Falcon 9 First Stage Landing Prediction

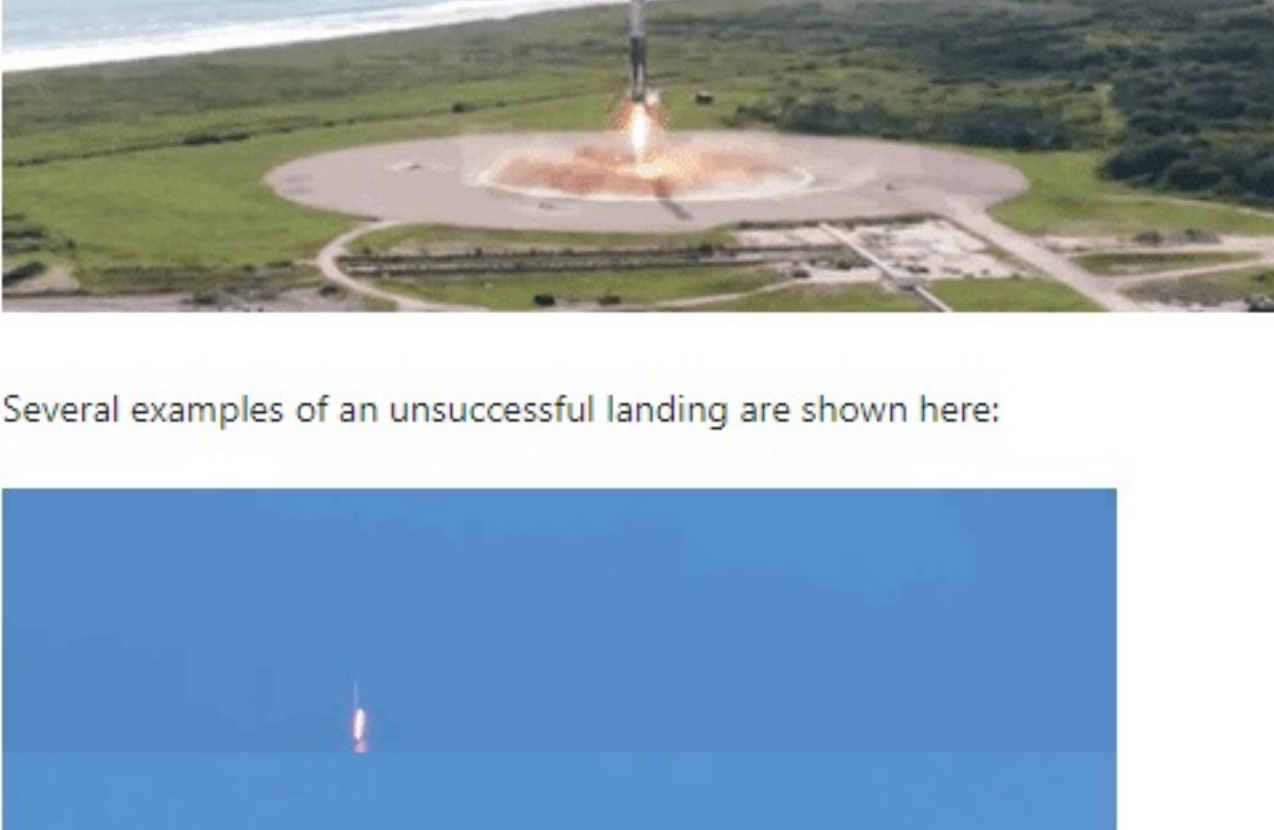
## Assignment: Exploring and Preparing Data

Estimated time needed: 70 minutes

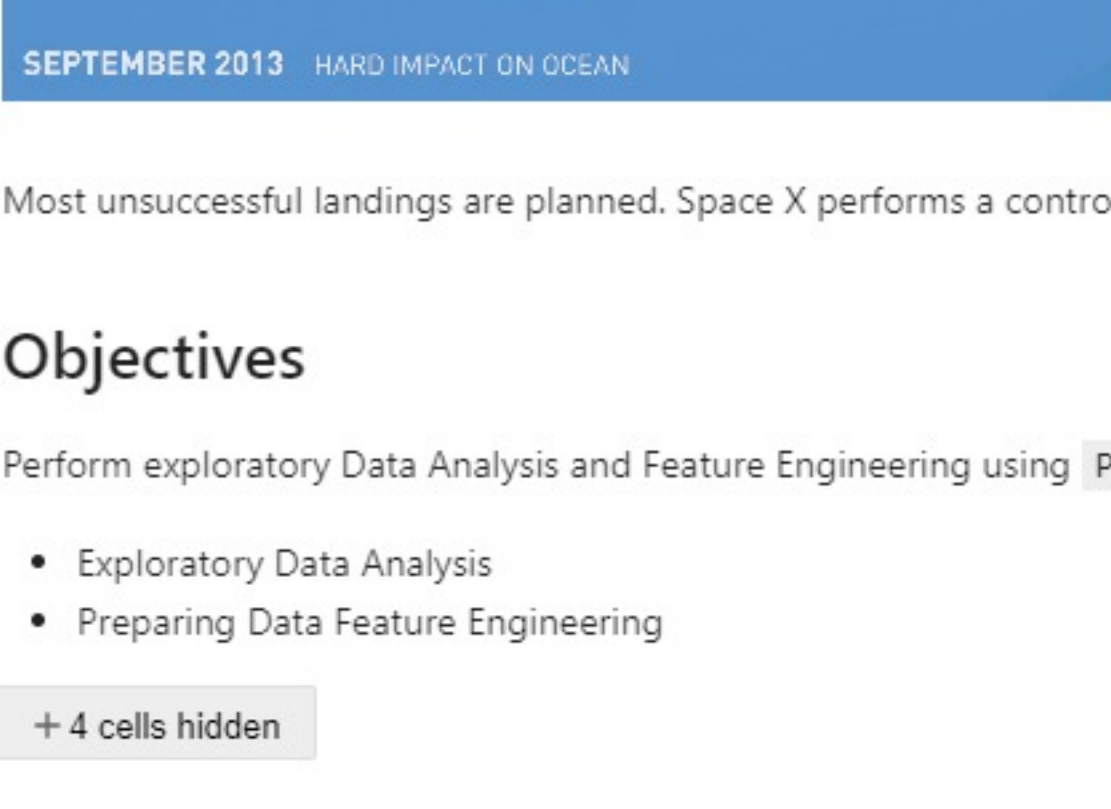
In this assignment, we will predict if the Falcon 9 first stage will land successfully. SpaceX advertises Falcon 9 rocket launches on its website with a cost of 62 million dollars; other providers cost upward of 165 million dollars each, much of the savings is due to the fact that SpaceX can reuse the 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

+ 4 cells hidden

## Exploratory Data Analysis

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

```
from js import fetch
```

```
[5]: FlightNumber Date BoosterVersion PayloadMass Orbit LaunchSite Outcome Flights GridFins Reused Legs LandingPad Block ReusedCount Serial Longitude
```

```
0 1 2010-06-04 Falcon 9 6104.959412 LEO CCAFS SLC 40 None None 1 False False False NaN 1.0 0 B0003 -80.57736
```

```
1 2 2012-05-22 Falcon 9 525.000000 LEO CCAFS SLC 40 None None 1 False False False NaN 1.0 0 B0005 -80.57736
```

```
2 3 2013-03-01 Falcon 9 677.000000 ISS CCAFS SLC 40 None None 1 False False False NaN 1.0 0 B0007 -80.57736
```

```
3 4 2013-09-29 Falcon 9 500.000000 PO VAFB SLC 4E False Ocean 1 False False False NaN 1.0 0 B1003 -120.61082
```

```
4 5 2013-12-03 Falcon 9 3170.000000 GTO CCAFS SLC 40 None None 1 False False False NaN 1.0 0 B1004 -80.57736
```

Next, 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.

```
[6]: 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()
```

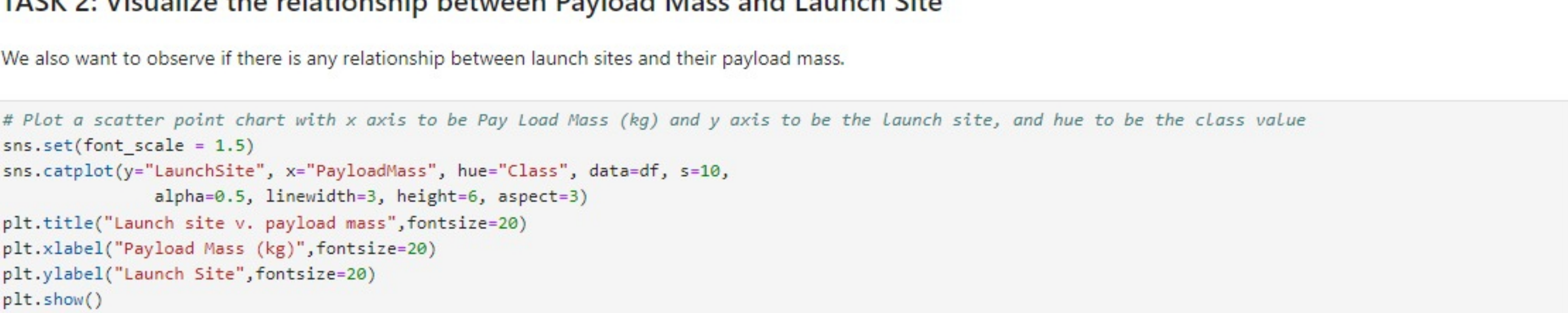


Next, let's drill down to each site visualize its detailed launch records.

## 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 `LaunchSite` and set the parameter `hue` to `'class'`

```
[7]: # 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.set(font_scale = 1.5)
sns.catplot(y="LaunchSite", x="FlightNumber", hue="Class", data=df, s=10,
            alpha=0.5, linewidth=3, height=6, aspect=3)
plt.title("Launch site v. Flight number", fontsize=20)
plt.xlabel("Flight Number", fontsize=20)
plt.ylabel("Launch Site", fontsize=20)
plt.show()
```

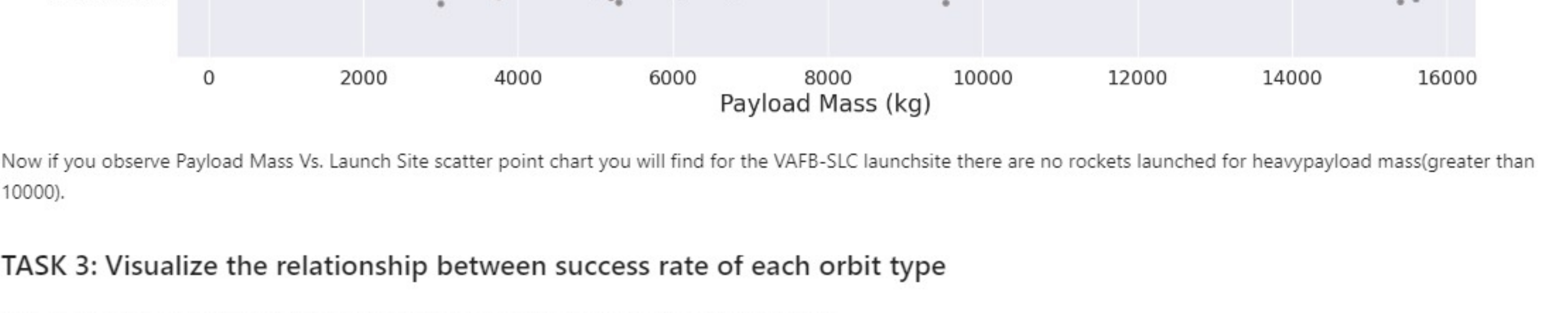


Now try to explain the patterns you found in the Flight Number vs. Launch Site scatter point plots.

## 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.

```
[8]: # 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 value
sns.set(font_scale = 1.5)
sns.catplot(y="LaunchSite", x="PayloadMass", hue="Class", data=df, s=10,
            alpha=0.5, linewidth=3, height=6, aspect=3)
plt.title("Launch site v. payload mass", fontsize=20)
plt.xlabel("Payload Mass (kg)", 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 heavy payload mass(greater than 10000).

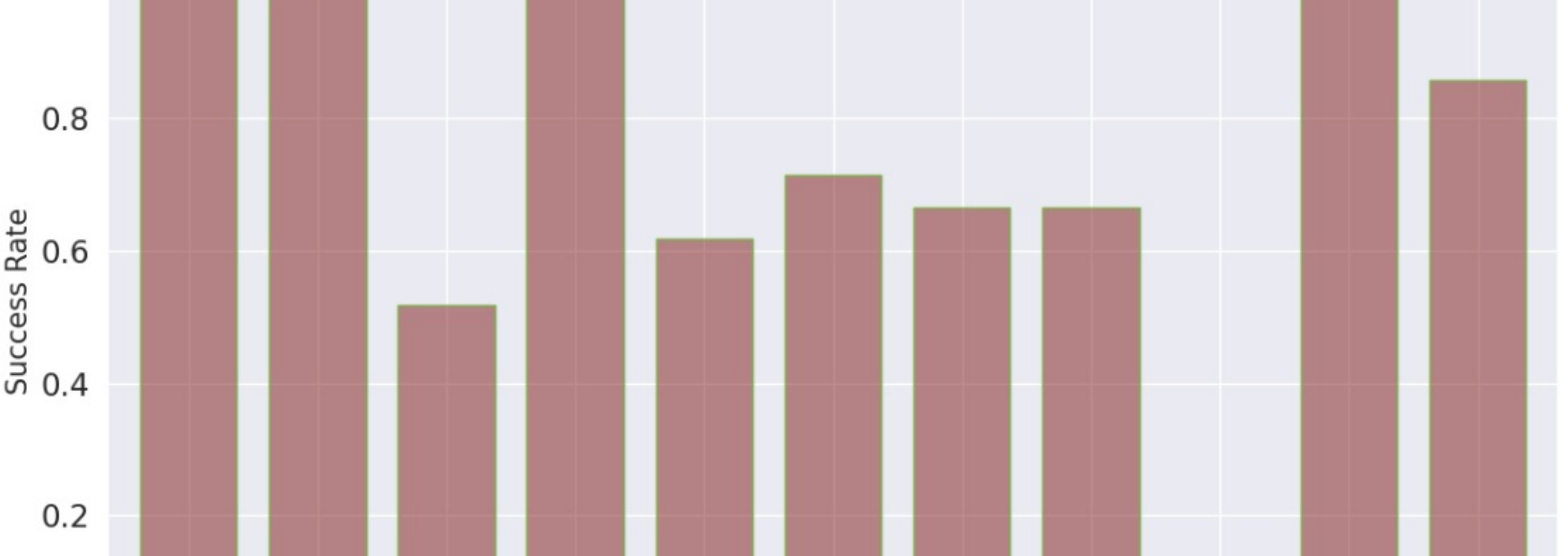
## TASK 3: Visualize the relationship between success rate of each orbit type

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 success rate of each orbit

```
[11]: # HINT use groupby method on Orbit column and get the mean of Class column
# Group by 'Orbit' and calculate the mean of 'Class', ignoring non-numeric columns
df.groupby(['Orbit']).mean(numeric_only=True)['Class'].plot(kind='bar', title="Success rate for each orbit",
                width=0.75, figsize=(12, 7),
                color=(0.5, 0.1, 0.1, 0.5),
                edgecolor='7eb54e')

plt.xlabel("Orbit", fontsize=15)
plt.ylabel("Success Rate", fontsize=15)
plt.tight_layout()
plt.show()
```

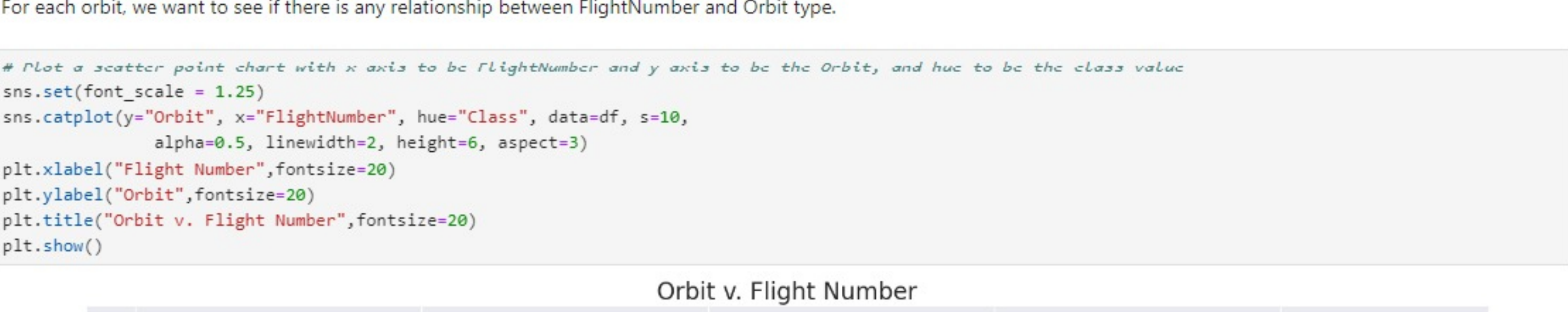


Analyze the plotted bar chart to identify which orbits have the highest success rates.

## 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.

```
[12]: # 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
sns.set(font_scale = 1.25)
sns.catplot(y="Orbit", x="FlightNumber", hue="Class", data=df, s=10,
            alpha=0.5, linewidth=2, height=6, aspect=3)
plt.xlabel("Flight Number", fontsize=20)
plt.ylabel("Orbit", fontsize=20)
plt.title("Orbit v. Flight Number", fontsize=20)
plt.show()
```

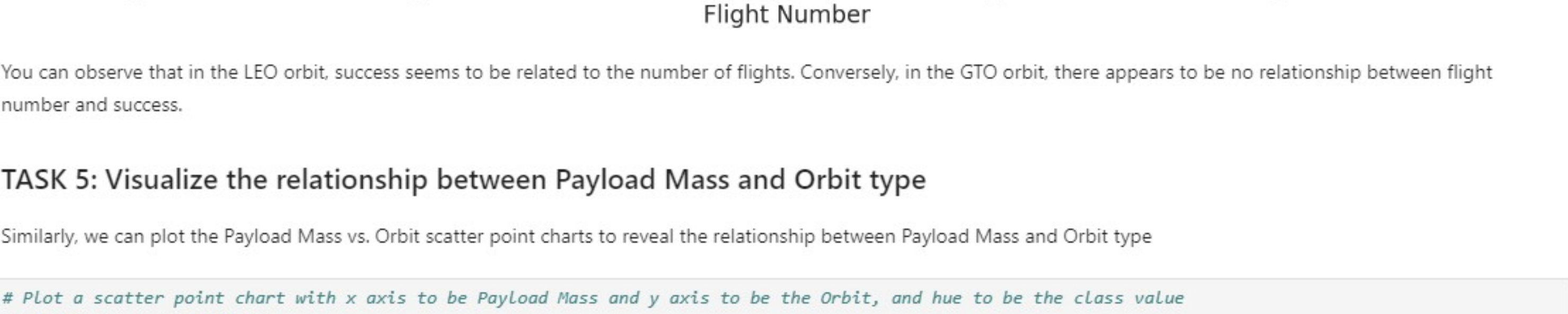


You can observe that in the LEO orbit, success seems to be related to the number of flights. Conversely, in the GTO orbit, there appears to be no relationship between flight number and success.

## 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

```
[13]: # 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
sns.set(font_scale = 1.25)
sns.catplot(y="Orbit", x="PayloadMass", hue="Class", data=df, s=10,
            alpha=0.5, linewidth=2, height=6, aspect=3)
#sns.catplot(y="Orbit", x="PayloadMass", hue="Class", data=df, aspect = 3)
plt.xlabel("Payload Mass(Kg)", fontsize=20)
plt.ylabel("Orbit", fontsize=20)
plt.show()
```



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

+ 5 cells hidden

## 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.

```
[17]: features = df[['FlightNumber', 'PayloadMass', 'Orbit', 'LaunchSite', 'Flights', 'GridFins', 'Reused', 'Legs', 'LandingPad', 'Block', 'ReusedCount', 'Serial']]
features.head()
```

```
[17]: 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.

```
[18]: # HINT: Use get_dummies() function on the categorical columns
features_one_hot = pd.get_dummies(features, columns=['Orbit', 'LaunchSite', 'LandingPad', 'Serial', 'GridFins', 'Reused', 'Legs'])
features_one_hot.head()
```

```
[18]: FlightNumber PayloadMass Flights Block ReusedCount Orbit_ES- L1 Orbit_GEO Orbit_GTO Orbit_HEO Orbit_ISS ... Serial_B1058 Serial_B1059 Serial_B1060 Serial_I
```

```
0 1 6104.959412 1 1.0 0 False False False False False ... False False False
```

```
1 2 525.000000 1 1.0 0 False False False False False ... False False False
```

```
2 3 677.000000 1 1.0 0 False False False False True ... False False False
```

```
3 4 500.000000 1 1.0 0 False False False False False ... False False False
```

```
4 5 3170.000000 1 1.0 0 False False True False False ... False False False
```

5 rows × 83 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`

```
[19]: # HINT: use astype function
features_one_hot.astype('float64').dtypes
```

```
[19]: FlightNumber float64
PayloadMass float64
Flights float64
Block float64
ReusedCount float64
GridFins_True float64
Reused_False float64
Reused_True float64
Legs_False float64
Legs_True float64
Length: 83, dtype: object
```

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