cognitiveclass.ai logo

# 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:

SEPTEMBER 2013    HARD IMPACT ON OCEAN

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

```python
In [1]:   # andas is a software library writ
          import pandas as pd
          #NumPy is a library for the Python
          import numpy as np
          # Matplotlib is a plotting library
          import matplotlib.pyplot as plt
          #Seaborn is a Python data visualiz
          import seaborn as sns
```

# Exploratory Data Analysis

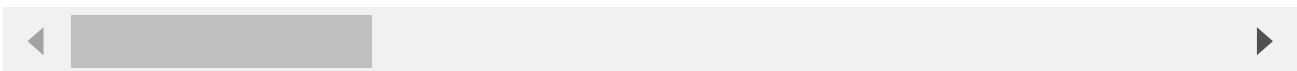First, let's read the SpaceX dataset into a Pandas dataframe and print its summary

```python
In [2]:   df=pd.read_csv("https://cf-courses
```

```
# If you were unable to complete t

# df = pd.read_csv('https://cf-cou

df.head(5)
```

In [2]:

Out[2]:

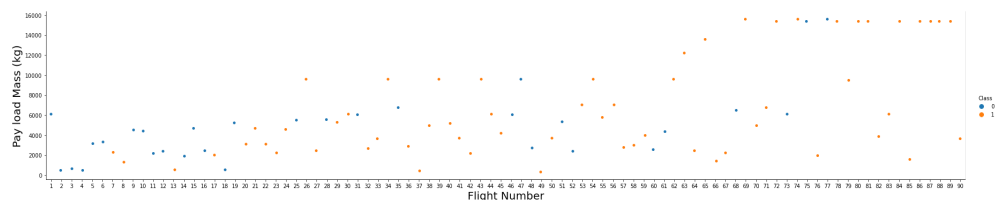| | FlightNumber | Date | BoosterVersion |
|---|---|---|---|
| **0** | 1 | 2010-06-04 | Falcon 9 |
| **1** | 2 | 2012-05-22 | Falcon 9 |
| **2** | 3 | 2013-03-01 | Falcon 9 |
| **3** | 4 | 2013-09-29 | Falcon 9 |
| **4** | 5 | 2013-12-03 | Falcon 9 |

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 is also important; it seems the more massive the payload, the less likely the first stage will return.

In [3]:
```python
sns.catplot(y="PayloadMass", x="Fl
plt.xlabel("Flight Number",fontsiz
plt.ylabel("Pay load Mass (kg)",fc
plt.show()
```

We see that different launch sites have different success rates. `CCAFS LC-40` , has a success rate of 60 %, while `KSC LC-39A` and `VAFB SLC 4E` has a success rate of 77%.
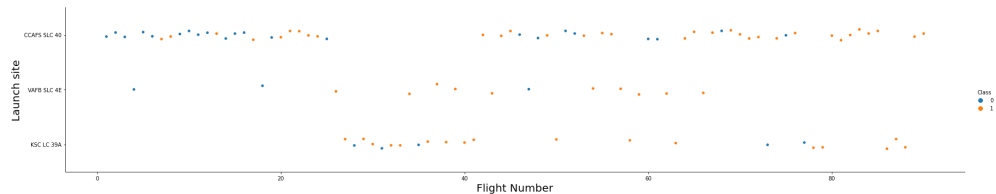
Next, let's drill down to each site visualize its detailed launch records.

# 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'`

In [5]:   `# Plot a scatter point chart with`

```python
sns.catplot(y="LaunchSite", x="Fli
plt.xlabel("Flight Number",fontsiz
plt.ylabel("Launch site",fontsize=
plt.show()
```
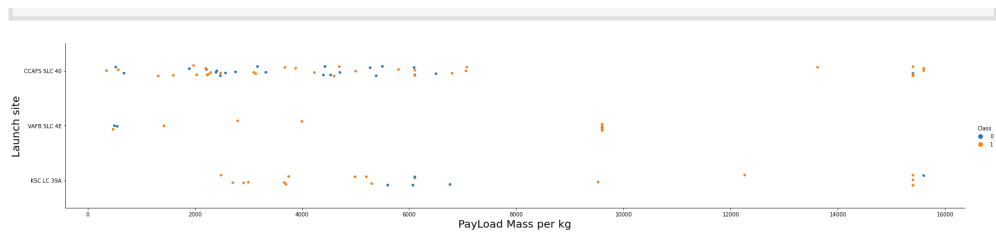


Now try to explain the patterns you found in the Flight Number vs. Launch Site scatter point plots.

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

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

```python
In [6]:   # Plot a scatter point chart with
          sns.catplot(y="LaunchSite", x="Pay
          plt.xlabel("PayLoad Mass per kg",f
          plt.ylabel("Launch site",fontsize=
          plt.show()
```

Now if you observe Payload 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).
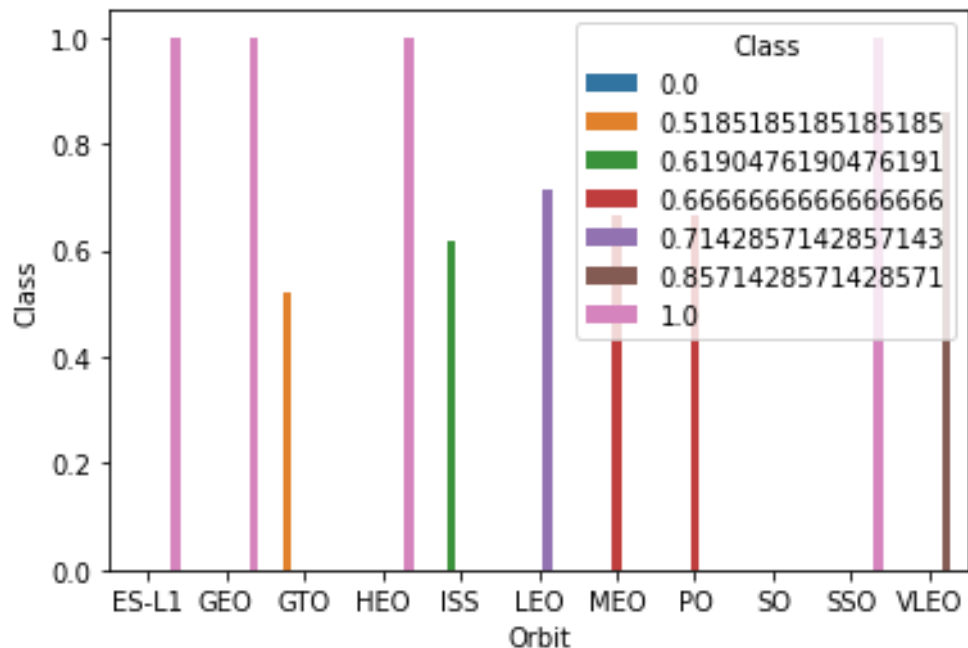
# TASK 3: Visualize the relationship between success rate of each orbit type

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 [7]: 
```
# HINT use groupby method on Orbit
orbit_success = df.groupby('Orbit'
orbit_success.reset_index(inplace=
sns.barplot(x="Orbit",y="Class",da
```

Out[7]:
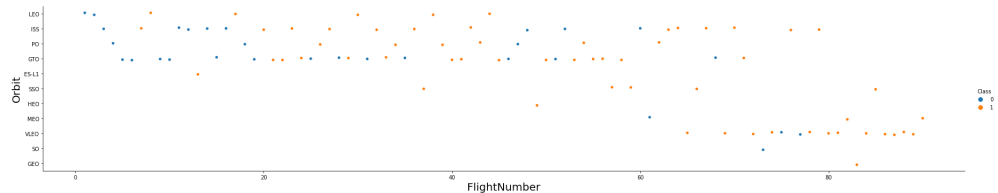```
<AxesSubplot:xlabel='Orbit', ylab
el='Class'>
```



Analyze the ploted bar chart try to find which orbits have high sucess rate.

# 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 [10]:
```python
# Plot a scatter point chart with
sns.catplot(y="Orbit", x="FlightNu
plt.xlabel("FlightNumber",fontsize
plt.ylabel("Orbit",fontsize=20)
plt.show()
```
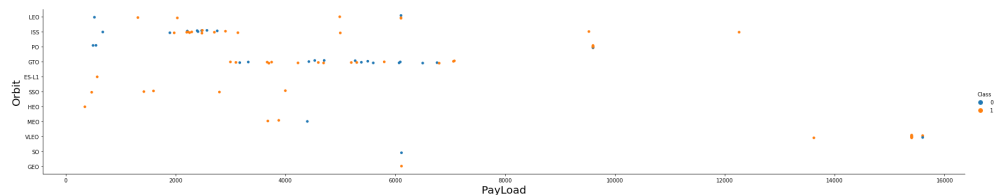


You should see that in the LEO orbit the Success appears related to the number of flights; on the other hand, there seems to be no relationship between flight number when in GTO orbit.

# TASK 5: Visualize the relationship between Payload and Orbit type

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

In [11]:
```python
# Plot a scatter point chart with
sns.catplot(y="Orbit", x="PayloadM
plt.xlabel("PayLoad",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 we cannot
distinguish this well as both
positive landing rate and negative
landing(unsuccessful mission) are
both there here.

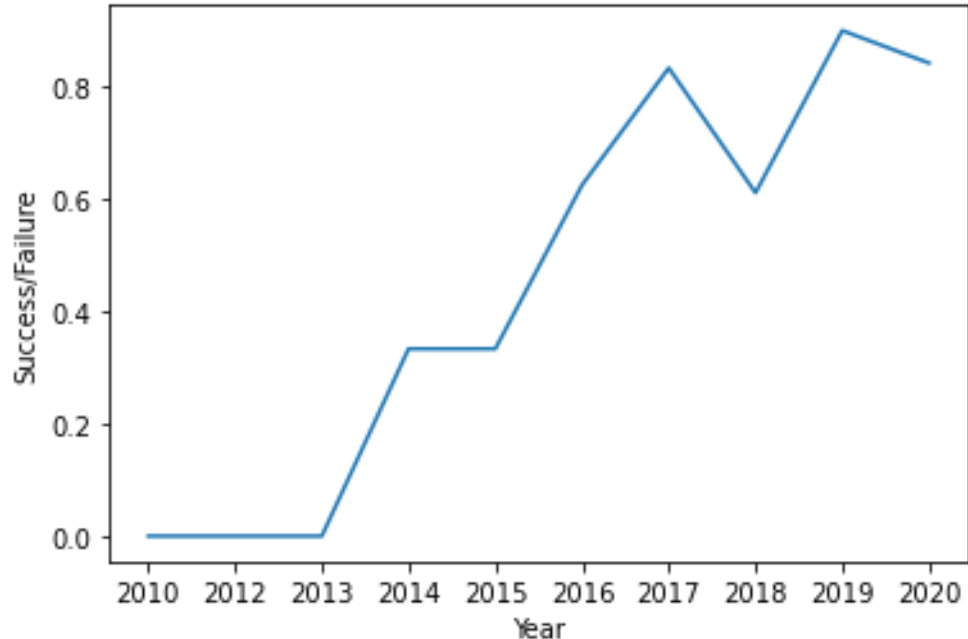# TASK 6: Visualize the launch success yearly trend

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 [52]:
```python
# A function to Extract years from
year=[]
def Extract_year(date):
    for i in df["Date"]:
        year.append(i.split("-")[0
    return year
Extract_year(1)
df["Year"]=year
average_by_year = df.groupby(by="Y
average_by_year.reset_index(inplac
```

In [54]:
```python
# Plot a line chart with x axis to
plt.plot(average_by_year["Year"],a
plt.xlabel("Year")
```

```python
plt.ylabel("Success/Failure")
plt.show()
```



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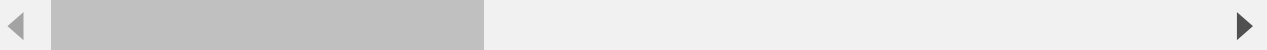
# Features Engineering

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 [55]:   features = df[['FlightNumber', 'Pa
           features.head()
```

Out[55]:

| | FlightNumber | PayloadMass | Orbit |
|---|---|---|---|
| **0** | 1 | 6104.959412 | LEO |
| **1** | 2 | 525.000000 | LEO |
| **2** | 3 | 677.000000 | ISS |
| **3** | 4 | 500.000000 | PO |
| **4** | 5 | 3170.000000 | GTO |

# 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.
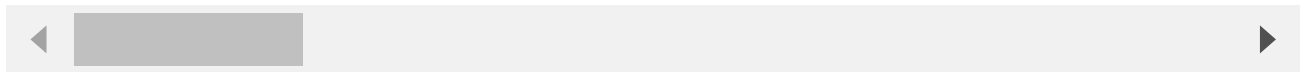
In [58]:
```python
# HINT: Use get_dummies() function
features_one_hot= pd.get_dummies(f
features_one_hot.head()
```

Out[58]:

| | FlightNumber | PayloadMass | Flights |
|---|---|---|---|
| **0** | 1 | 6104.959412 | 1 |
| **1** | 2 | 525.000000 | 1 |
| **2** | 3 | 677.000000 | 1 |
| **3** | 4 | 500.000000 | 1 |
| **4** | 5 | 3170.000000 | 1 |

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`
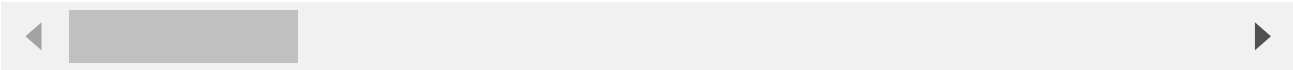
In [59]:
```
# HINT: use astype function
features_one_hot =features_one_hot
```

features_one_hot

Out[59]:

| | FlightNumber | PayloadMass | Flights |
|---|---|---|---|
| **0** | 1.0 | 6104.959412 | 1.0 |
| **1** | 2.0 | 525.000000 | 1.0 |
| **2** | 3.0 | 677.000000 | 1.0 |
| **3** | 4.0 | 500.000000 | 1.0 |
| **4** | 5.0 | 3170.000000 | 1.0 |
| **...** | ... | ... | ... |
| **85** | 86.0 | 15400.000000 | 2.0 |
| **86** | 87.0 | 15400.000000 | 3.0 |
| **87** | 88.0 | 15400.000000 | 6.0 |
| **88** | 89.0 | 15400.000000 | 3.0 |
| **89** | 90.0 | 3681.000000 | 1.0 |

90 rows × 80 columns

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('datase
 index=False)
```

# Authors

Joseph Santarcangelo has a PhD in
Electrical Engineering, his research
focused on using machine learning,
signal processing, and computer
vision to determine how videos
impact human cognition. Joseph
has been working for IBM since he
completed his PhD.

Nayef Abou Tayoun is a Data
Scientist at IBM and pursuing a

Master of Management in Artificial intelligence degree at Queen's University.

# Change Log

| Date (YYYY-MM-DD) | Version | Changed By | Change Description |
|---|---|---|---|
| 2021-10-12 | 1.1 | Lakshmi Holla | Modified markdown |
| 2020-09-20 | 1.0 | Joseph | Modified Multiple Areas |
| 2020-11-10 | 1.1 | Nayef | updating the input data |