# 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

```python
# 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 mat
import numpy as np
# Matplotlib is a plotting library for python and pyplot gives us a MatLab like plotting framework. We will use th
import matplotlib.pyplot as plt
#Seaborn is a Python data visualization library based on matplotlib. It provides a high-level interface for drawing
import seaborn as sns
```

```python
## Exploratory Data Analysis

```

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

```
In [10]:    1  URL = "https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBM-DS0321EN-SkillsNetwork/datasets/data
            2  df=pd.read_csv(URL)
            3  df.head(5)
```

Out[10]:

| | FlightNumber | Date | BoosterVersion | PayloadMass | Orbit | LaunchSite | Outcome | Flights | GridFins | Reused | Legs | LandingPad | Block | ReusedC |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 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 | |

```
In [11]:    1  df.describe()
```

Out[11]:

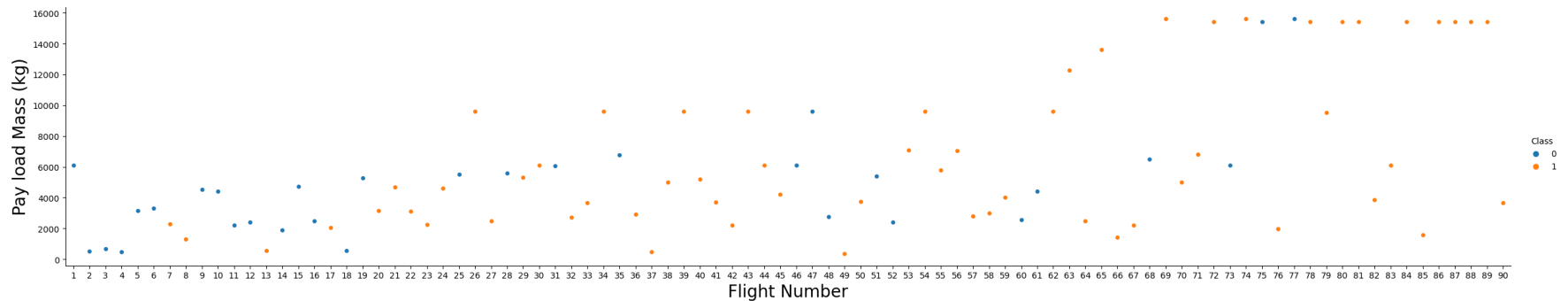| | FlightNumber | PayloadMass | Flights | Block | ReusedCount | Longitude | Latitude | Class |
|---|---|---|---|---|---|---|---|---|
| count | 90.000000 | 90.000000 | 90.000000 | 90.000000 | 90.000000 | 90.000000 | 90.000000 | 90.000000 |
| mean | 45.500000 | 6104.959412 | 1.788889 | 3.500000 | 1.655556 | -86.366477 | 29.449963 | 0.666667 |
| std | 26.124701 | 4694.671720 | 1.213172 | 1.595288 | 1.710254 | 14.149518 | 2.141306 | 0.474045 |
| min | 1.000000 | 350.000000 | 1.000000 | 1.000000 | 0.000000 | -120.610829 | 28.561857 | 0.000000 |
| 25% | 23.250000 | 2510.750000 | 1.000000 | 2.000000 | 0.000000 | -80.603956 | 28.561857 | 0.000000 |
| 50% | 45.500000 | 4701.500000 | 1.000000 | 4.000000 | 1.000000 | -80.577366 | 28.561857 | 1.000000 |
| 75% | 67.750000 | 8912.750000 | 2.000000 | 5.000000 | 3.000000 | -80.577366 | 28.608058 | 1.000000 |
| max | 90.000000 | 15600.000000 | 6.000000 | 5.000000 | 5.000000 | -80.577366 | 34.632093 | 1.000000 |

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

In [12]: ▶|
```python
1  sns.catplot(y="PayloadMass", x="FlightNumber", hue="Class", data=df, aspect = 5)
2  plt.xlabel("Flight Number",fontsize=20)
3  plt.ylabel("Pay load Mass (kg)",fontsize=20)
4  plt.show()
```

C:\Users\parichea\AppData\Local\anaconda3\Lib\site-packages\seaborn\axisgrid.py:118: UserWarning: The figure layout h
as changed to tight
  self._figure.tight_layout(*args, **kwargs)



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.

In [16]: ▶|
```python
1  ### TASK 1: Visualize the relationship between Flight Number and Launch Site
2
```
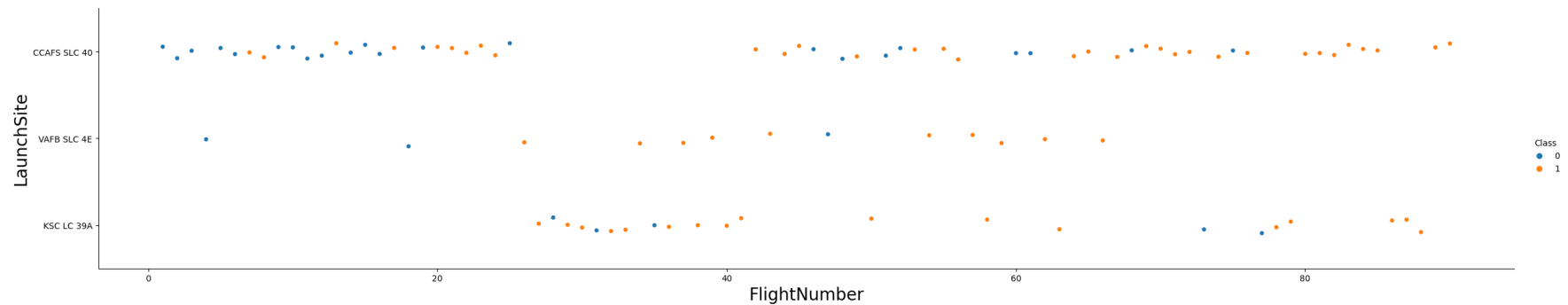
Use the function `catplot` to plot `FlightNumber` vs `LaunchSite` , set the parameter `x` parameter to `FlightNumber` ,set the `y` to `Launch Site` and set the parameter `hue` to `'class'`

In [15]: ▶|
```python
1  # 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 c
2  sns.catplot(y="LaunchSite", x="FlightNumber", hue="Class", data=df, aspect = 5)
3  plt.xlabel("FlightNumber",fontsize=20)
4  plt.ylabel("LaunchSite",fontsize=20)
5  plt.show()
```

C:\Users\parichea\AppData\Local\anaconda3\Lib\site-packages\seaborn\axisgrid.py:118: UserWarning: The figure layout h
as changed to tight
  self._figure.tight_layout(*args, **kwargs)



Now try to explain the patterns you found in the Flight Number vs. Launch Site scatter point plots.
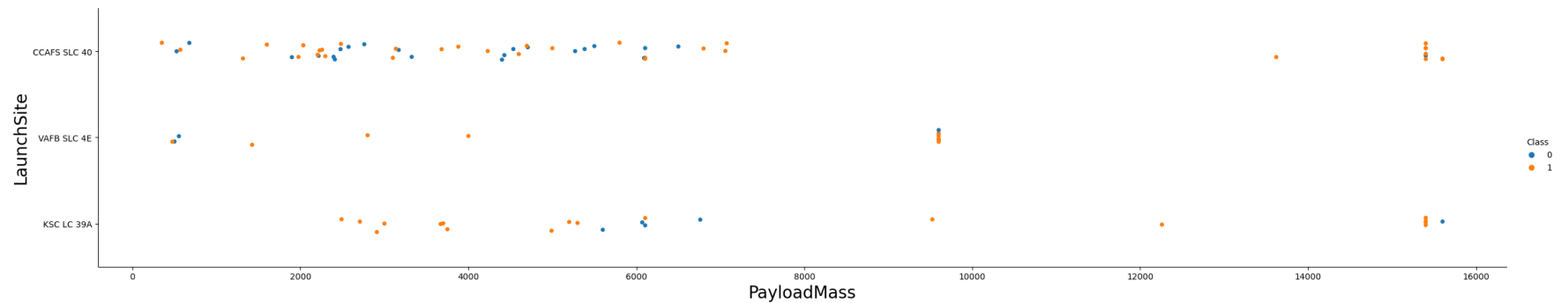
In [18]: ▶|
```python
1  ### TASK 2: Visualize the relationship between Payload and Launch Site
2
3
```

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

```
1  # 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
2  sns.catplot(y="LaunchSite", x="PayloadMass", hue="Class", data=df, aspect = 5)
3  plt.xlabel("PayloadMass",fontsize=20)
4  plt.ylabel("LaunchSite",fontsize=20)
5  plt.show()
```

```
C:\Users\parichea\AppData\Local\anaconda3\Lib\site-packages\seaborn\axisgrid.py:118: UserWarning: The figure layout h
as changed to tight
  self._figure.tight_layout(*args, **kwargs)
```



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).

```
1  ### TASK  3: Visualize the relationship between success rate of each orbit type
2  sns.countplot(data=df, x='Orbit', hue='Class')
3  plt.show()
```

```
In [26]:    1 sns.countplot(data=df_success, x='Orbit', hue='Class')
            2 plt.show()
```

```
In [27]:  ▶   1  sns.countplot(data=df_fail, x='Orbit', hue='Class')
              2  plt.show()
```



```
In [22]:  ▶   1  df_success=df[df['Class']==1]
              2  df_fail= df[df['Class']==0]
```

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

```
In [23]:  ▶   1  y=set(df_success['Orbit'])
              2  y
```

Out[23]: {'ES-L1', 'GEO', 'GTO', 'HEO', 'ISS', 'LEO', 'MEO', 'PO', 'SSO', 'VLEO'}

In [24]: ▶ 
```
1  x=set(df_fail['Orbit'])
2  x
```

Out[24]: {'GTO', 'ISS', 'LEO', 'MEO', 'PO', 'SO', 'VLEO'}

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

In [25]: ▶ 
```
1  per=(df_success['Orbit'].value_counts())
2  per
```

Out[25]: 
```
Orbit
GTO     14
ISS     13
VLEO    12
PO       6
LEO      5
SSO      5
MEO      2
ES-L1    1
HEO      1
GEO      1
Name: count, dtype: int64
```

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

In [ ]: ▶ 
```
1  ### TASK  4: Visualize the relationship between FlightNumber and Orbit type
2
```

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

```
1  # Plot a scatter point chart with x axis to be FlightNumber and y axis to be the Orbit, and hue to be the class va
2  sns.catplot(y='Orbit', x='FlightNumber', hue='Class',data=df,aspect = 5)
3  plt.xlabel('FlightNumber', fontsize=20)
4  plt.xlabel('Orbit', fontsize=20)
5  plt.show()
6
```

C:\Users\parichea\AppData\Local\anaconda3\Lib\site-packages\seaborn\axisgrid.py:118: UserWarning: The figure layout h
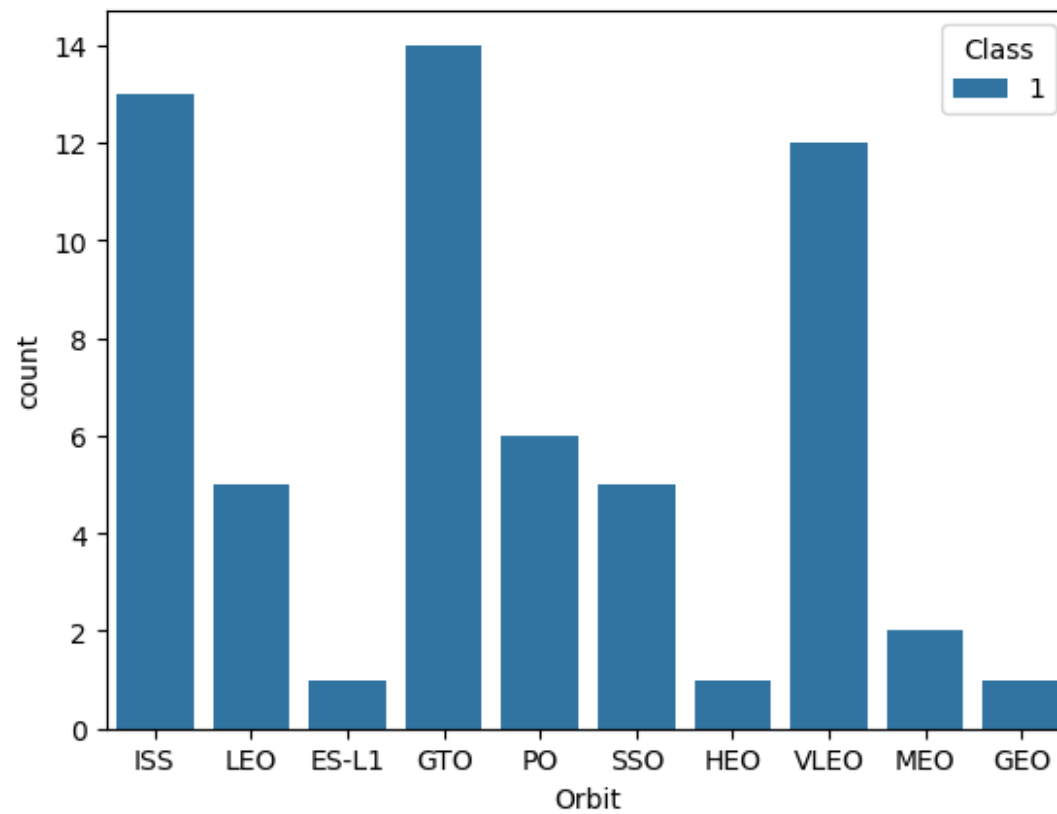as changed to tight
  self._figure.tight_layout(*args, **kwargs)

```
1  sns.catplot(y='Orbit', x='FlightNumber', hue='LaunchSite',data=df,aspect = 5)
2  plt.xlabel('FlightNumber', fontsize=20)
3  plt.xlabel('Orbit', fontsize=20)
4  plt.show()
```

C:\Users\parichea\AppData\Local\anaconda3\Lib\site-packages\seaborn\axisgrid.py:118: UserWarning: The figure layout h
as changed to tight
  self._figure.tight_layout(*args, **kwargs)

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.

In [ ]: ▶| 
```
1  ### TASK  5: Visualize the relationship between Payload and Orbit type
2
```

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

In [30]: ▶|
```
1  # Plot a scatter point chart with x axis to be Payload and y axis to be the Orbit, and hue to be the class value
2  sns.catplot(y='Orbit', x='PayloadMass', hue='Class',data=df,aspect = 5)
3  plt.xlabel('PayloadMass', fontsize=20)
4  plt.xlabel('Orbit', fontsize=20)
5  plt.show()
```

C:\Users\parichea\AppData\Local\anaconda3\Lib\site-packages\seaborn\axisgrid.py:118: UserWarning: The figure layout h
as changed to tight
  self._figure.tight_layout(*args, **kwargs)

In [31]: ▶|
```
1  sns.catplot(y='Orbit', x='PayloadMass', hue='LaunchSite',data=df,aspect = 5)
2  plt.xlabel('PayloadMass', fontsize=20)
3  plt.xlabel('Orbit', fontsize=20)
4  plt.show()
5
```

C:\Users\parichea\AppData\Local\anaconda3\Lib\site-packages\seaborn\axisgrid.py:118: UserWarning: The figure layout h
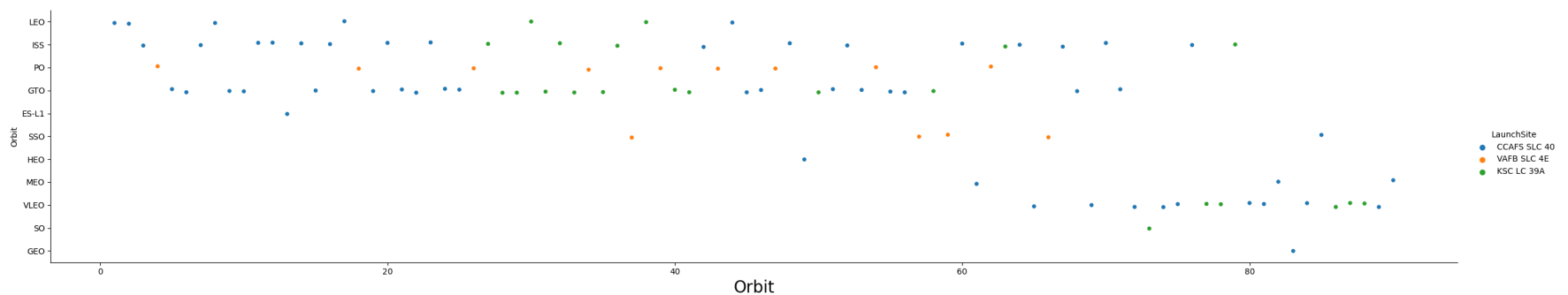as changed to tight
  self._figure.tight_layout(*args, **kwargs)



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.

In [ ]: ▶|
```
1  ### TASK  6: Visualize the launch success yearly trend
2
```

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 [32]:  ▶|    1  # A function to Extract years from the date
               2  year=[]
               3  def Extract_year():
               4      for i in df["Date"]:
               5          year.append(i.split("-")[0])
               6      return year
               7  Extract_year()
               8  df['Date'] = year
               9  df.head()
              10
```

Out[32]:

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

```
In [33]:  ▶|    1  df_success=df[df['Class']==1]
```

```
1 df_line=df_success[['Date', 'Class']]
2 df_line
```

Out[37]:

| | Date | Class |
|---|---|---|
| 6 | 2014 | 1 |
| 7 | 2014 | 1 |
| 12 | 2015 | 1 |
| 16 | 2015 | 1 |
| 19 | 2016 | 1 |
| 20 | 2016 | 1 |
| 21 | 2016 | 1 |
| 22 | 2016 | 1 |
| 23 | 2016 | 1 |
| 25 | 2017 | 1 |
| 26 | 2017 | 1 |
| 28 | 2017 | 1 |
| 29 | 2017 | 1 |
| 31 | 2017 | 1 |
| 32 | 2017 | 1 |
| 33 | 2017 | 1 |
| 35 | 2017 | 1 |
| 36 | 2017 | 1 |
| 37 | 2017 | 1 |
| 38 | 2017 | 1 |
| 39 | 2017 | 1 |
| 40 | 2017 | 1 |
| 41 | 2017 | 1 |
| 42 | 2017 | 1 |
| 43 | 2018 | 1 |
| 44 | 2018 | 1 |

|    | Date | Class |
|----|------|-------|
| 48 | 2018 | 1 |
| 49 | 2018 | 1 |
| 52 | 2018 | 1 |
| 53 | 2018 | 1 |
| 54 | 2018 | 1 |
| 55 | 2018 | 1 |
| 56 | 2018 | 1 |
| 57 | 2018 | 1 |
| 58 | 2018 | 1 |
| 61 | 2019 | 1 |
| 62 | 2019 | 1 |
| 63 | 2019 | 1 |
| 64 | 2019 | 1 |
| 65 | 2019 | 1 |
| 66 | 2019 | 1 |
| 68 | 2019 | 1 |
| 69 | 2019 | 1 |
| 70 | 2019 | 1 |
| 71 | 2020 | 1 |
| 73 | 2020 | 1 |
| 75 | 2020 | 1 |
| 77 | 2020 | 1 |
| 78 | 2020 | 1 |
| 79 | 2020 | 1 |
| 80 | 2020 | 1 |
| 81 | 2020 | 1 |
| 82 | 2020 | 1 |

|    | Date | Class |
|----|------|-------|
| 83 | 2020 | 1 |
| 84 | 2020 | 1 |
| 85 | 2020 | 1 |
| 86 | 2020 | 1 |
| 87 | 2020 | 1 |
| 88 | 2020 | 1 |
| 89 | 2020 | 1 |

```
In [35]:    1  df_success['Class'].count()
```

Out[35]: 60

```
In [ ]:    1
```

```
1  # Plot a line chart with x axis to be the extracted year and y axis to be the success rate
2  sns.countplot(x='Date', data=df_success)
3  plt.show()
```



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

In [ ]:

```
1  ## Features Engineering
2
```

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 [39]:  ▶|    1  features = df[['FlightNumber', 'PayloadMass', 'Orbit', 'LaunchSite', 'Flights', 'GridFins', 'Reused', 'Legs', 'Lan
                2  features.head()
```

Out[39]:

| | 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 |

```
In [ ]:  ▶|    1  ### TASK  7: Create dummy variables to categorical columns
              2
```

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 [46]:  ▶|    1  # HINT: Use get_dummies() function on the categorical columns
                2  features_one_hot= pd.get_dummies(df[['Orbit', 'LaunchSite','LandingPad', 'Serial']])
                3  features_one_hot.head()
```

Out[46]:

| | Orbit_ES-L1 | Orbit_GEO | Orbit_GTO | Orbit_HEO | Orbit_ISS | Orbit_LEO | Orbit_MEO | Orbit_PO | Orbit_SO | Orbit_SSO | ... | Serial_B1048 | Serial_B1049 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | False | False | False | False | False | True | False | False | False | False | ... | False | False |
| 1 | False | False | False | False | False | True | False | False | False | False | ... | False | False |
| 2 | False | False | False | False | True | False | False | False | False | False | ... | False | False |
| 3 | False | False | False | False | False | False | False | True | False | False | ... | False | False |
| 4 | False | False | True | False | False | False | False | False | False | False | ... | False | False |

5 rows × 72 columns

```
In [47]: ▶| 1 df_Dummy= features_one_hot.astype(float)
```

```
In [60]: ▶| 1 df_Dummy
```

Out[60]:

| | Orbit_ES-L1 | Orbit_GEO | Orbit_GTO | Orbit_HEO | Orbit_ISS | Orbit_LEO | Orbit_MEO | Orbit_PO | Orbit_SO | Orbit_SSO | ... | Serial_B1048 | Serial_B104$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 |
| 1 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 |
| 2 | 0.0 | 0.0 | 0.0 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 |
| 3 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 1.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 |
| 4 | 0.0 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | .. |
| 85 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 |
| 86 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 |
| 87 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 |
| 88 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 |
| 89 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 |

90 rows × 72 columns

```
In [ ]: ▶| 1 df=df.drop(['Orbit', 'LaunchSite','LandingPad', 'Serial'], axis=1)
```

```
In [53]:  ▶  1  df
```

Out[53]:

| | FlightNumber | Date | BoosterVersion | PayloadMass | Outcome | Flights | GridFins | Reused | Legs | Block | ReusedCount | Longitude | Latitude | C |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 1 | 2010 | Falcon 9 | 6104.959412 | None None | 1 | False | False | False | 1.0 | 0 | -80.577366 | 28.561857 | |
| **1** | 2 | 2012 | Falcon 9 | 525.000000 | None None | 1 | False | False | False | 1.0 | 0 | -80.577366 | 28.561857 | |
| **2** | 3 | 2013 | Falcon 9 | 677.000000 | None None | 1 | False | False | False | 1.0 | 0 | -80.577366 | 28.561857 | |
| **3** | 4 | 2013 | Falcon 9 | 500.000000 | False Ocean | 1 | False | False | False | 1.0 | 0 | -120.610829 | 34.632093 | |
| **4** | 5 | 2013 | Falcon 9 | 3170.000000 | None None | 1 | False | False | False | 1.0 | 0 | -80.577366 | 28.561857 | |
| **...** | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| **85** | 86 | 2020 | Falcon 9 | 15400.000000 | True ASDS | 2 | True | True | True | 5.0 | 2 | -80.603956 | 28.608058 | |
| **86** | 87 | 2020 | Falcon 9 | 15400.000000 | True ASDS | 3 | True | True | True | 5.0 | 2 | -80.603956 | 28.608058 | |
| **87** | 88 | 2020 | Falcon 9 | 15400.000000 | True ASDS | 6 | True | True | True | 5.0 | 5 | -80.603956 | 28.608058 | |
| **88** | 89 | 2020 | Falcon 9 | 15400.000000 | True ASDS | 3 | True | True | True | 5.0 | 2 | -80.577366 | 28.561857 | |
| **89** | 90 | 2020 | Falcon 9 | 3681.000000 | True ASDS | 1 | True | False | True | 5.0 | 0 | -80.577366 | 28.561857 | |

90 rows × 14 columns

```
In [59]:  ▶  1  df=pd.concat([df, df_Dummy], axis=1)
```

```
In [56]:    1 df
```

Out[56]:

| | FlightNumber | Date | BoosterVersion | PayloadMass | Outcome | Flights | GridFins | Reused | Legs | Block | ... | Serial_B1048 | Serial_B1049 | Serial_B' |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 1 | 2010 | Falcon 9 | 6104.959412 | None None | 1 | False | False | False | 1.0 | ... | 0.0 | 0.0 | |
| **1** | 2 | 2012 | Falcon 9 | 525.000000 | None None | 1 | False | False | False | 1.0 | ... | 0.0 | 0.0 | |
| **2** | 3 | 2013 | Falcon 9 | 677.000000 | None None | 1 | False | False | False | 1.0 | ... | 0.0 | 0.0 | |
| **3** | 4 | 2013 | Falcon 9 | 500.000000 | False Ocean | 1 | False | False | False | 1.0 | ... | 0.0 | 0.0 | |
| **4** | 5 | 2013 | Falcon 9 | 3170.000000 | None None | 1 | False | False | False | 1.0 | ... | 0.0 | 0.0 | |
| **...** | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| **85** | 86 | 2020 | Falcon 9 | 15400.000000 | True ASDS | 2 | True | True | True | 5.0 | ... | 0.0 | 0.0 | |
| **86** | 87 | 2020 | Falcon 9 | 15400.000000 | True ASDS | 3 | True | True | True | 5.0 | ... | 0.0 | 0.0 | |
| **87** | 88 | 2020 | Falcon 9 | 15400.000000 | True ASDS | 6 | True | True | True | 5.0 | ... | 0.0 | 0.0 | |
| **88** | 89 | 2020 | Falcon 9 | 15400.000000 | True ASDS | 3 | True | True | True | 5.0 | ... | 0.0 | 0.0 | |
| **89** | 90 | 2020 | Falcon 9 | 3681.000000 | True ASDS | 1 | True | False | True | 5.0 | ... | 0.0 | 0.0 | |

90 rows × 86 columns

```
In [58]:    1 df.to_csv('week02_02.csv', index=False)
```

```
In [ ]:    1 ### TASK  8: Cast all numeric columns to `float64`
           2
```

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

```
In [65]:  ▶  1  # HINT: use astype function
             2  df_Dummy= features_one_hot.astype(float)
```

```
In [66]:  ▶  1  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 (https://www.linkedin.com/in/pratiksha-verma-6487561b1/?utm_medium=Exinfluencer&utm_source=Exinfluencer&utm_content=000026UJ&utm_term=10006555&utm_id=NA-SkillsNetwork-Channel-SkillsNetworkCoursesIBMDS0321ENSkillsNetwork865-2022-01-01)

# Change Log

| Date (YYYY-MM-DD) | Version | Changed By | Change Description |
|---|---|---|---|
| 2022-11-09 | 1.0 | Pratiksha Verma | Converted initial version to Jupyterlite |