



IBM Developer  
SKILLS NETWORK

# Winning Space Race with Data Science

Arian Moslehi  
12/27/2022



# Outline

---

- Executive Summary
- Introduction
- Methodology
- Results
- Conclusion
- Appendix

# Executive Summary

---

- Summary of methodologies

- Data collection
- Data wrangling
- EDA with data visualization
- EDA with SQL
- Predictive analysis (Classification)

- Summary of all results

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

# Introduction

---

- Project background and context

We predicted 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 because SpaceX can reuse the first stage. Therefore, if we can determine if the first stage will land, we can determine the cost of a launch. This information can be used if an alternate company wants to bid against SpaceX for a rocket launch.

- Problems you want to find answers

- What influences if the rocket will land successfully?
- The effect each relationship with certain rocket variables will impact in determining the success rate of a successful landing.
- What conditions does SpaceX have to achieve to get the best results and ensure the best rocket success landing rate?



Section 1

# Methodology

# Methodology

---

## Executive Summary

- Data collection methodology:
  - SpaceX Rest API
  - (Web Scrapping)
- Perform data wrangling
  - One Hot Encoding data fields for Machine Learning and dropping irrelevant columns
- Perform exploratory data analysis (EDA) using visualization and SQL
- Perform predictive analysis using classification models

# Data Collection

---

## **The following datasets was collected by:**

- We worked with SpaceX launch data that is gathered from the SpaceX REST API.
- This API will give us data about launches, including information about the rocket used, payload delivered, launch specifications, landing specifications, and landing outcome.
- Our goal is to use this data to predict whether SpaceX will attempt to land a rocket or not.
- The SpaceX REST API endpoints, or URL, starts with `api.spacexdata.com/v4/`.
- Another popular data source for obtaining Falcon 9 Launch data is web scraping Wikipedia using

# Data Collection – SpaceX API

---

[GitHub URL to Notebook](#)





# Data Collection - Scraping

*simplified flow chart*

## 1 .Getting Response from HTML

```
page = requests.get(static_url)
```

## 2. Creating BeautifulSoup Object

```
soup = BeautifulSoup(page.text, 'html.parser')
```

## 3. Finding tables

```
html_tables = soup.find_all('table')
```

## 4. Getting column names

```
column_names = []
temp = soup.find_all('th')
for x in range(len(temp)):
    try:
        name = extract_column_from_header(temp[x])
        if (name is not None and len(name) > 0):
            column_names.append(name)
    except:
        pass
```

## 5. Creation of dictionary

```
launch_dict= dict.fromkeys(column_names)

# Remove an irrelevant column
del launch_dict['Date and time ( )']

launch_dict['Flight No.'] = []
launch_dict['Launch site'] = []
launch_dict['Payload'] = []
launch_dict['Payload mass'] = []
launch_dict['Orbit'] = []
launch_dict['Customer'] = []
launch_dict['Launch outcome'] = []
launch_dict['Version Booster']=[[]]
launch_dict['Booster landing']=[[]]
launch_dict['Date']=[[]]
launch_dict['Time']=[[]]
```

## 6. Appending data to keys (refer) to notebook block 12

```
In [12]: extracted_row = 0
#Extract each table
for table_number,table in enumerate(
    # get table row
    for rows in table.find_all("tr")
    #check to see if first table
```

## 7. Converting dictionary to dataframe

```
df = pd.DataFrame.from_dict(launch_dict)
```

## 8. Dataframe to .CSV

```
df.to_csv('spacex_web_scraped.csv', index=False)
```

# Data Wrangling

---

## Introduction

- In the data set, there are several different cases where the booster did not land successfully. Sometimes a landing was attempted but failed due to an accident; for example, True Ocean means the mission outcome was successfully landed to a specific region of the ocean while False Ocean means the mission outcome was unsuccessfully landed to a specific region of the ocean. True RTLS means the mission outcome was successfully landed to a ground pad False RTLS means the mission outcome was unsuccessfully landed to a ground pad. True ASDS means the mission outcome was successfully landed on a drone ship False ASDS means the mission outcome was unsuccessfully landed on a drone ship. We mainly convert those outcomes into Training Labels with 1 means the booster successfully landed 0 means it was unsuccessful

[GitHub URL](#)

## Process

Perform Exploratory Data Analysis EDA on dataset

Calculate the number of launches at each site

Calculate the number and occurrence of each orbit

Create a landing outcome label from Outcome column

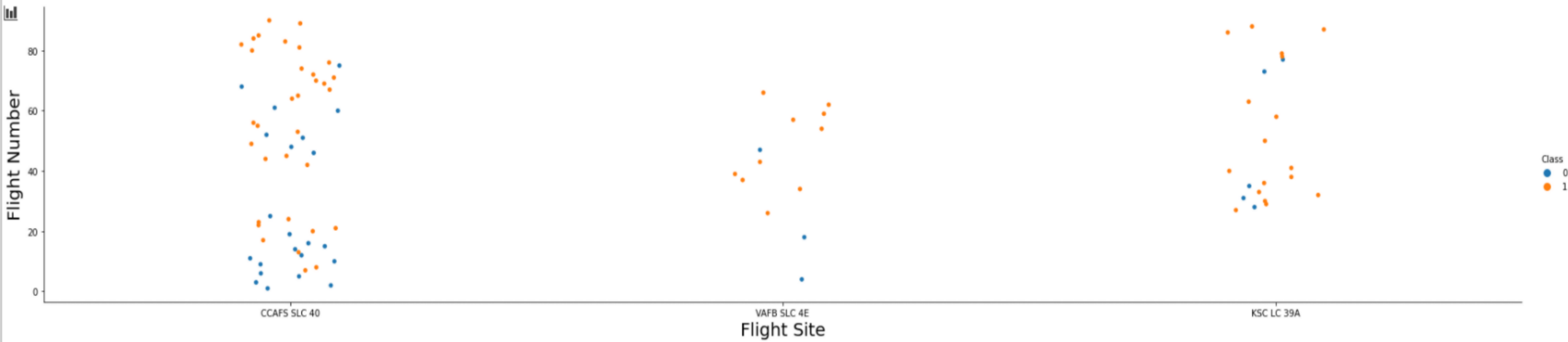
Work out success rate for every landing in dataset

# EDA with Data Visualization

[GitHub URL](#)

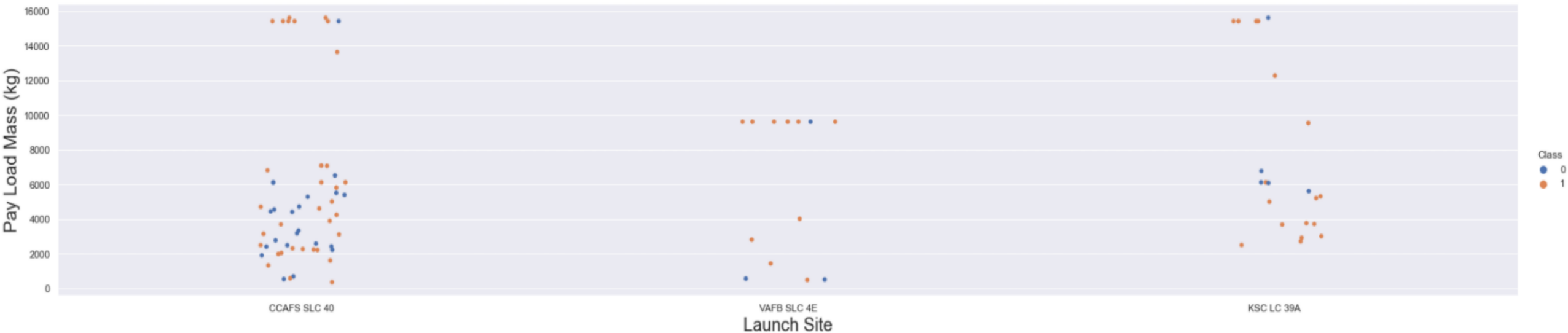


# Flight Number vs. Flight Site



The more amount of flights at a launch site the greater the success rate at a launch site.

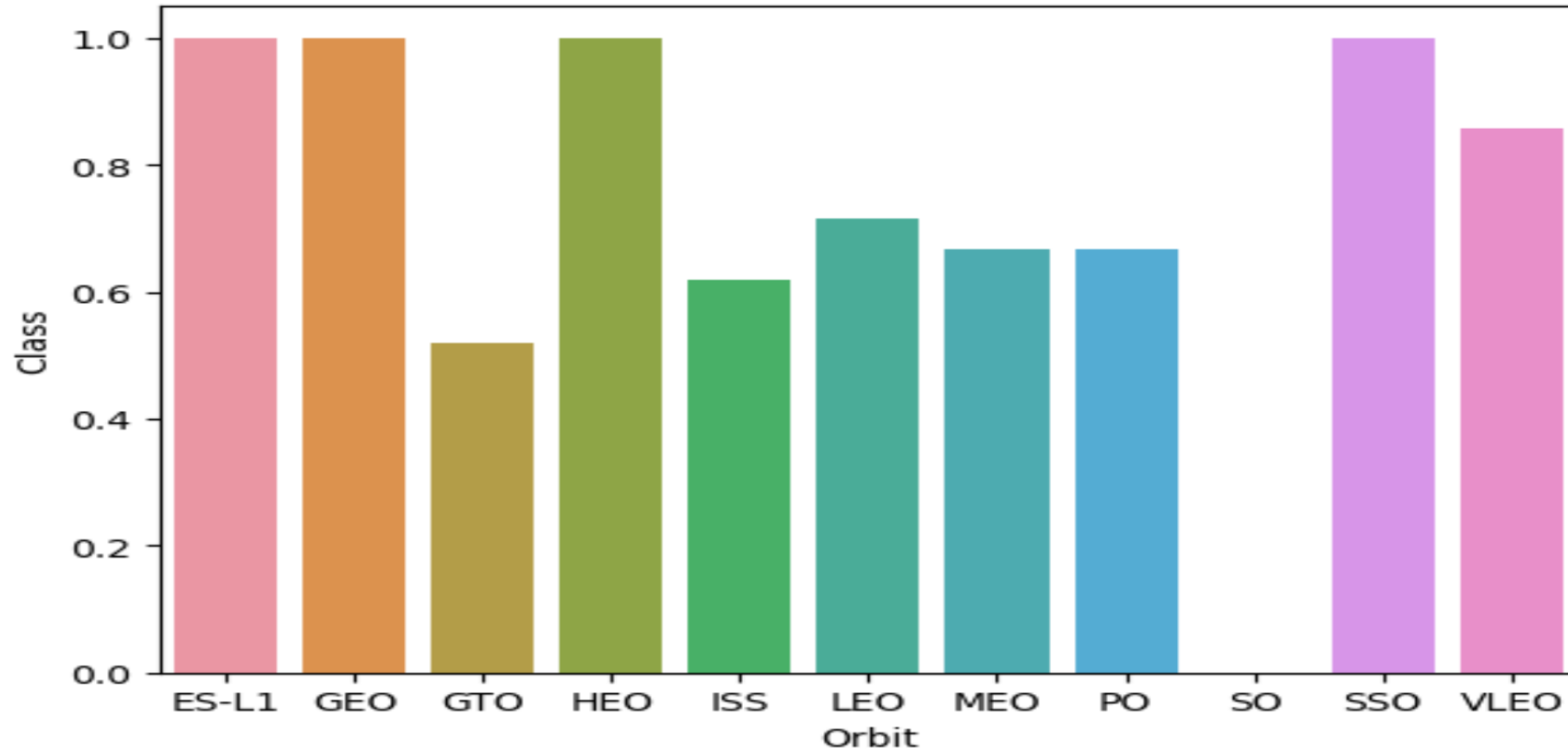
# Payload Mass vs. Launch Site



The greater the payload mass for Launch Site CCAFS SLC 40 the higher the success rate for the Rocket. There is not quite a clear pattern to be found using this visualization to make a decision if the Launch Site is dependant on Pay Load Mass for a success launch.



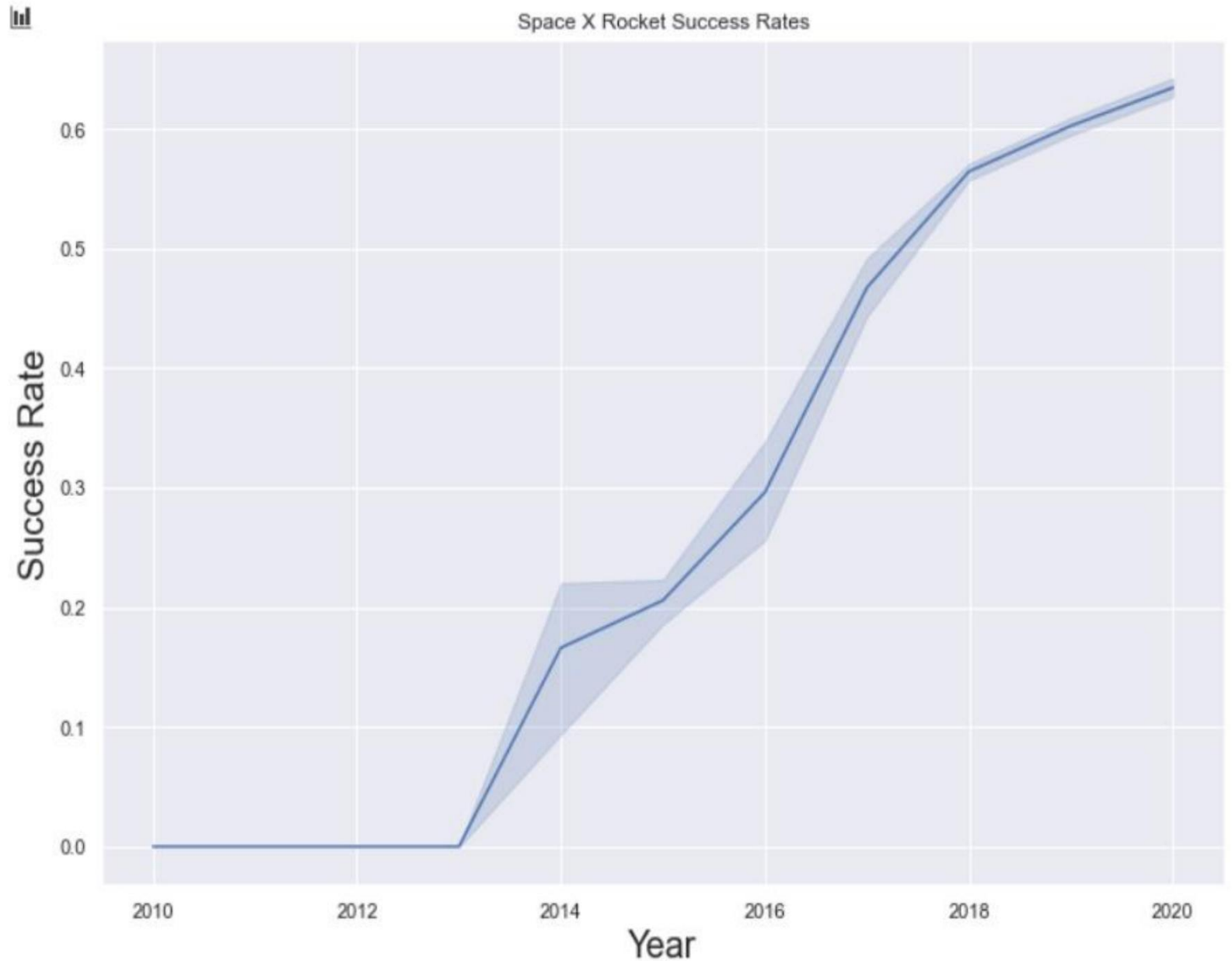
# Success Rate



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

# Launch success yearly trend

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



# EDA with SQL

---

**EDA WITH  
.SQL**

[GitHub URL](#)

# summary

- Finding distinct launch sites from spaceX table
- Using the function SUM summates the total in the column PAYLOAD\_MASS\_KG\_ for Customer NASA (CRS)
- Using the function AVG works out the average in the column PAYLOAD\_MASS\_KG\_ for Booster\_version F9 v1.1
- Using the function MIN to find first successful landing
- Finding boosters that carried maximum payload
- Rank success count between 2010-06-04 and 2017-03-20

# Predictive Analysis (Classification)





# Classification Accuracy

- The decision tree model had the highest accuracy in training and testing data.

```
In [22]: tree_cv = GridSearchCV(tree, param_grid=parameters,scoring='accuracy', cv=10)
tree_cv.fit(X_train, Y_train)
tree_cv.best_params_
```

```
Out[22]: {'criterion': 'gini',
'max_depth': 4,
'max_features': 'auto',
'min_samples_leaf': 2,
'min_samples_split': 5,
'splitter': 'random'}
```

```
In [23]: print("tuned hpyerparameters :(best parameters) ",tree_cv.best_params_)
print("accuracy :",tree_cv.best_score_)
```

```
tuned hpyerparameters :(best parameters) {'criterion': 'gini', 'max_depth': 4, 'max_features': 'auto', 'min_samples_leaf': 2, 'min_samples_split': 5,
'splitter': 'random'}
accuracy : 0.8875
```

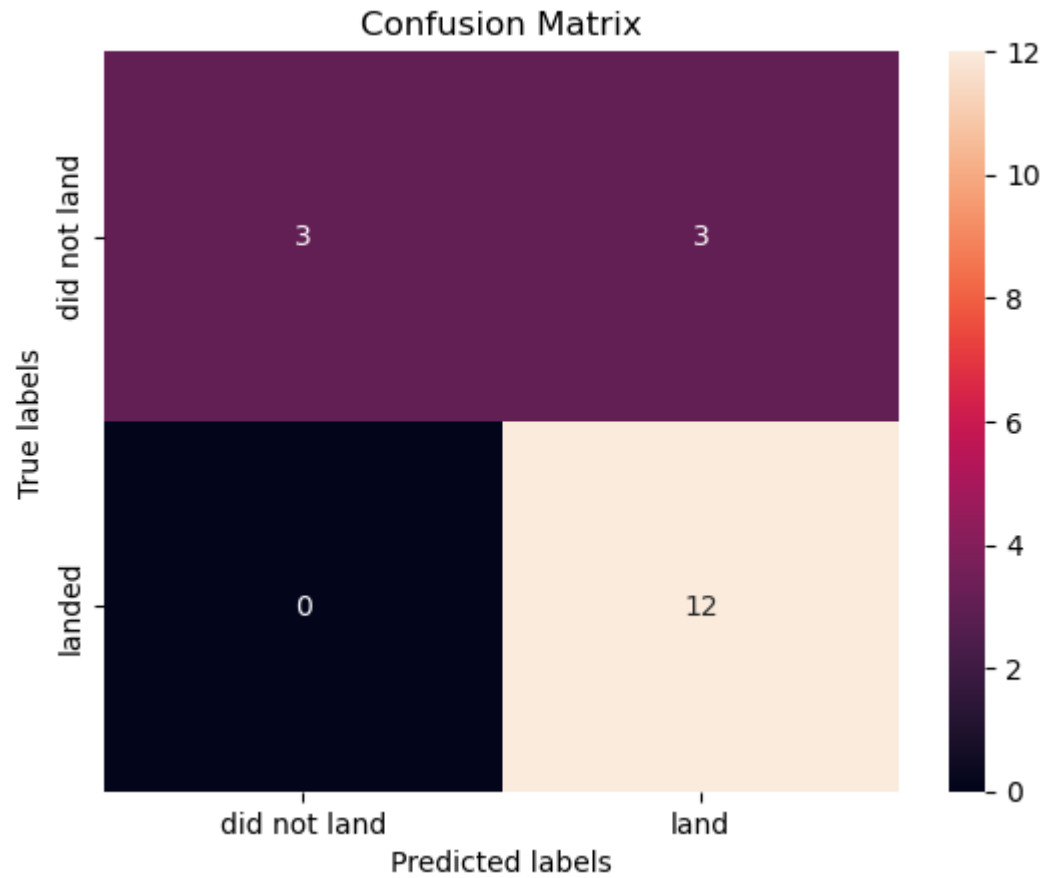
## TASK 9

Calculate the accuracy of tree\_cv on the test data using the method `score` :

```
In [24]: tree_cv.score(X_test, Y_test)
```

```
Out[24]: 0.8333333333333334
```

# Decision Tree Confusion Matrix



[GitHub URL](#)

# Conclusion

---

- The Tree Classifier Algorithm is the best for Machine Learning for this dataset
- Low weighted payloads perform better than the heavier payloads
- The success rates for SpaceX launches is directly proportional time in years they will eventually perfect the launches
- We can see that KSC LC-39A had the most successful launches from all the sites
- Orbit GEO,HEO,SSO,ES-L1 has the best Success Rate

Thank you!

