



IBM Developer  
SKILLS NETWORK

# Winning Space Race with Data Science

Bruno L. S. Bedo, PhD

December 07<sup>th</sup>, 2022



# Outline

---

**01**

## **Executive Summary**

Overview

**02**

## **Introduction**

Main question

**03**

## **Methodology**

Data | EDA | DashBoard |  
Classification

**04**

## **Results**

Type of analysis

**05**

## **Conclusion**

Discussion and main findings

**06**

## **Appendix**

Github and notebook links

# Executive Summary

---



## Methodologies

- Data Collection from ExpaceEx API and web scraping
- Exploratory Data Analysis (EDA)
- Prediction: Machine Learning Tools



## Main results

Exploratory data analysis

Interactive analysis

Predictive analysis

# Introduction

---

**SpaceX designs, manufactures and launches advanced rockets and spacecraft. The company was founded in 2002 to revolutionize space technology**

**Using public information, as well machine learning tools, we are predicted the best place to make launches**



Section 1

# Methodology

# Methodology

---

- Data Collection.
  - SpaceX API ([link](#))
  - Web Scrapping from Wikipedia ([link](#))
- Data wrangling;
  - Filtering the data;
  - Dealing with missing values.
  - To label de data label based on outcome data after summarizing and analyzing features;
- Performed exploratory data analysis (EDA) using visualization and SQL;
- Performed interactive visual analytics using Folium and Plotly Dash;
- Performed predictive analysis using classification models.



Section 1

# Methodology



# Data Collection - SpaceX API

i) To get response response from API and convert the result to *.json* file

```
1 static_json_url='https://cf-courses-data.s3.us.cloud-object-storage.appdomain.clo
2 response.status_code
3 # Use json_normalize meethod to convert the json result into a dataframe
4 response = requests.get(static_json_url).json()
5 df = pd.json_normalize(response)
```

Python

ii) Cleaning data

```
1 # Call getLaunchSite
2 getLaunchSite(data)
```

[38] ✓ 50.6s

```
1 # Call getPayloadData
2 getPayloadData(data)
```

[39] ✓ 48.6s

```
1 # Call getCoreData
2 getCoreData(data)
```

[40] ✓ 48.5s



SpaceX API: <https://api.spacexdata.com/v4/rockets/>





# Data Collection - SpaceX API

## iii) Converting list to a dataframe

Finally lets construct our dataset using the data we have obtained. We we combine the columns into a dictionary.

[+ Code](#) [+ Markdown](#)

```
1 launch_dict = {'FlightNumber': list(data['flight_number']),
2               'Date': list(data['date']),
3               'BoosterVersion':BoosterVersion,
4               'PayloadMass':PayloadMass,
5               'Orbit':Orbit,
6               'LaunchSite':LaunchSite,
7               'Outcome':Outcome,
8               'Flights':Flights,
9               'GridFins':GridFins,
10              'Reused':Reused,
11              'Legs':Legs,
12              'LandingPad':LandingPad,
13              'Block':Block,
14              'ReusedCount':ReusedCount,
15              'Serial':Serial,
16              'Longitude': Longitude,
17              'Latitude': Latitude}
18
```

## iv) Filtering the dataframe and converting it to a .csv file

```
1 data_falcon9.to_csv('dataset_part_1.csv', index=False)
```

✓ 0.1s



# Data Collection - Scrapping

## i) Get a response from HTML adress

```
1 # use requests.get() method with the provided static
2 # assign the response to a object
3 page = requests.get(static_url)
4 page.status_code
```

[12] ✓ 1.2s Python

... 200

## iii) To find the tables of interest

```
1 # Use the find_all function in the BeautifulSoup obj
2 # Assign the result to a list called `html_tables`
3 html_tables = soup.find_all('table')
```

[15] ✓ 0.9s Python

## ii) To create a BeautifulSoup Object

Create a BeautifulSoup object from the HTML response

```
1 # Use BeautifulSoup() to create a BeautifulSoup object
2 soup = BeautifulSoup(page.text, 'html.parser')
```

[13] ✓ 1.3s Python

## iv) To get the columns info

```
1 column_names = []
2
3 # Apply find_all() function with `th` element on first table
4 # Iterate each th element and apply the provided extract_column_from_header function
5 # Append the Non-empty column name (if name is not None)
6 temp = soup.find_all('th')
7 for x in range(len(temp)):
8     try:
9         name = extract_column_from_header(temp[x])
10        if (name is not None and len(name) > 0):
11            column_names.append(name)
12    except:
13        pass
```

✓ 0.1s Python



# Data Collection - Scrapping

## v) Creating a dictionary

```
1 launch_dict= dict.fromkeys(column_names)
2
3 # Remove an irrelevant column
4 del launch_dict['Date and time ( )']
5
6 # Let's initial the launch_dict with each value to be
7 launch_dict['Flight No.'] = []
8 launch_dict['Launch site'] = []
9 launch_dict['Payload'] = []
10 launch_dict['Payload mass'] = []
11 launch_dict['Orbit'] = []
12 launch_dict['Customer'] = []
13 launch_dict['Launch outcome'] = []
14 # Added some new columns
15 launch_dict['Version Booster']=[]
16 launch_dict['Booster landing']=[]
17 launch_dict['Date']=[]
18 launch_dict['Time']=[]
```

✓ 0.7s

Python

## ii) Appending data from all keys, converting to a dataframe and saving in .csv

```
1 extracted_row = 0
2 #Extract each table
3 for table_number,table in enumerate(soup.find_all('table',"wikitable plainrowhead
4 # get table row
5 for rows in table.find_all("tr"):
6 #check to see if first table heading is as number corresponding to launch
7 if rows.th:
8 if rows.th.string:
9 flight_number=rows.th.string.strip()
10 flag=flight_number.isdigit()
11 else:
12 flag=False
13 #get table element
14 row=rows.find_all('td')
15 #if it is number save cells in a dictionary
16 if flag:
```

```
1 headings = []
2 for key,values in dict(launch_dict).items():
3 if key not in headings:
4 headings.append(key)
5 if values is None:
6 del launch_dict[key]
7
8 def pad_dict_list(dict_list, padel):
9 lmax = 0
10 for lname in dict_list.keys():
11 lmax = max(lmax, len(dict_list[lname]))
12 for lname in dict_list.keys():
13 ll = len(dict_list[lname])
14 if ll < lmax:
15 dict_list[lname] += [padel] * (lmax - ll)
16 return dict_list
17
18 pad_dict_list(launch_dict,0)
19
20 df = pd.DataFrame.from_dict(launch_dict)
21 df.head()
22 df.to_csv('spacex_web_scraped.csv', index=False)
```

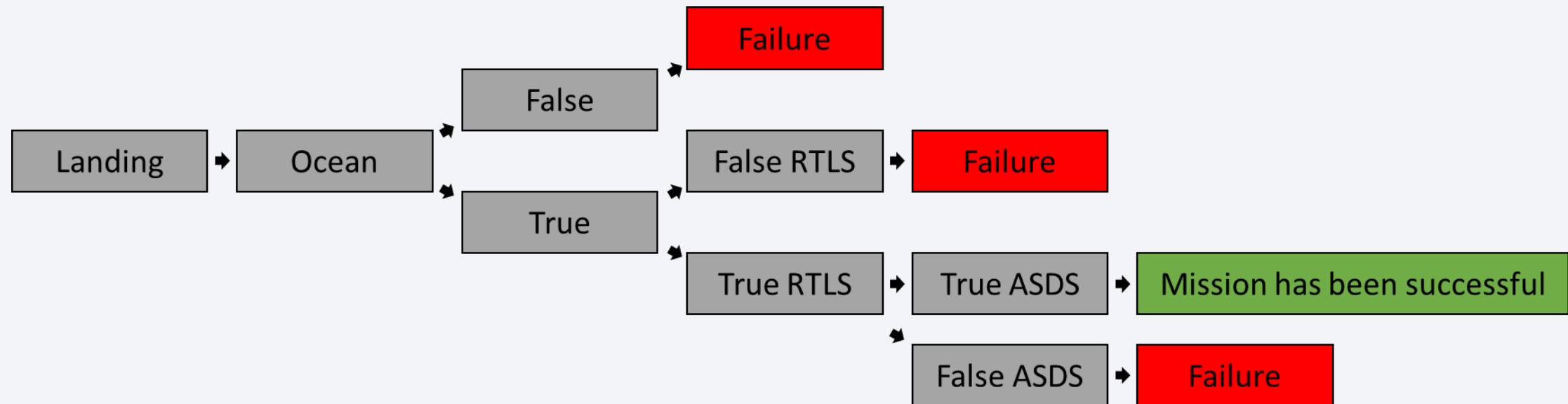
✓ 0.1s





# Data Wrangling

Regarding the data analysis process, there are several cases in which the booster did not land successfully.





# Data Wrangling

## i) Number of launches

```
[34]: 1 # Apply value_counts() on column LaunchSite
      2 df["LaunchSite"].value_counts()

...   CCAFS SLC 40    55
      KSC LC 39A    22
      VAFB SLC 4E    13
      Name: LaunchSite, dtype: int64
```

## ii) Number accourence of Earth orbit

```
[35]: 1 # Apply value_counts on Orbit column
      2 df["Orbit"].value_counts("Orbit")

...   GTO    0.300000
      ISS    0.233333
      VLEO   0.155556
      PO     0.100000
      LEO    0.077778
      SSO    0.055556
      MEO    0.033333
      HEO    0.011111
      SO     0.011111
      ES-L1  0.011111
      GEO    0.011111
      Name: Orbit, dtype: float64
```

## iii) Number missoun outcome

```
[44]: 1 # landing_outcomes = values on Outcome column
      2 landing_outcomes = df["Outcome"].value_counts()
      3 landing_outcomes

...   True ASDS    41
      None None    19
      True RTLS    14
      False ASDS    6
      True Ocean    5
      None ASDS     2
      False Ocean   2
      False RTLS    1
      Name: Outcome, dtype: int64
```

## iv) Creating a outcome label

```
[37]: 1 for i,outcome in enumerate(landing_outcomes.keys()):
      2     print(i,outcome)

...   0 True ASDS
      1 None None
      2 True RTLS
      3 False ASDS
      4 True Ocean
      5 None ASDS
      6 False Ocean
      7 False RTLS
```

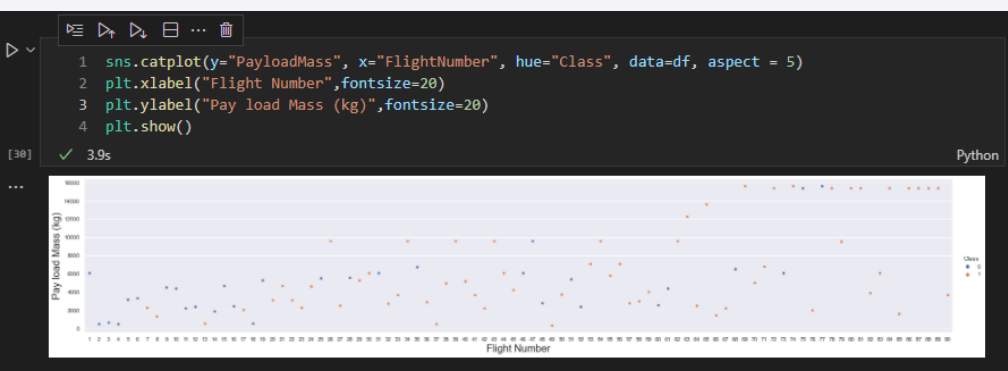


# EDA with Data Visualization

## Scatter plots

Scatter plots showed the relationship between two variables (correlation)

Flight Number VS. Payload Mass  
Flight Number VS. Launch Site  
Payload VS. Launch Site  
Orbit VS. Flight Number  
Payload VS. Orbit Type  
Orbit VS. Payload Mass



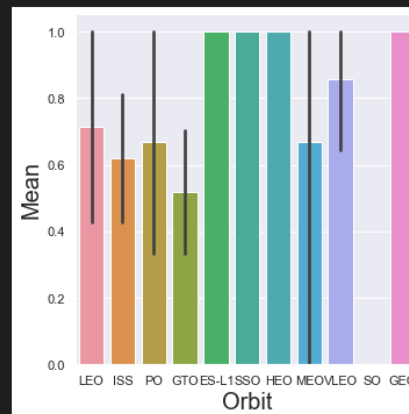
## Bar plots

Mean VS. Orbit

A bar diagram makes it easy to compare sets of data between different groups.

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

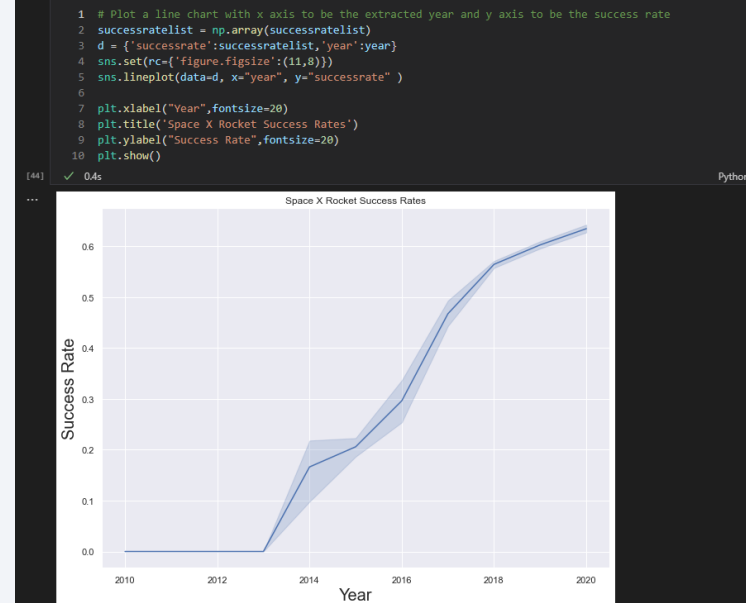
```
1 sns.catplot(x="Orbit",y="Class", kind="bar",data=df)
2 plt.xlabel("Orbit",fontsize=20)
3 plt.ylabel("Mean",fontsize=20)
4 plt.show()
```



## Line plots

Success Rate VS. Year

Line graphs are useful in that they show data variables and trends during the time







# EDA with SQL

---

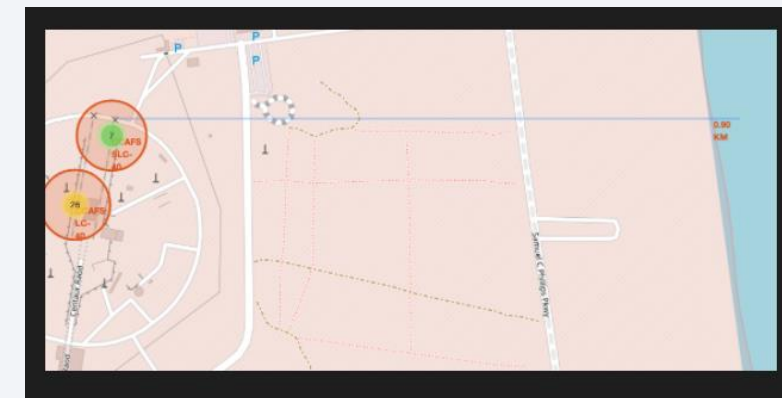
Please see below the SQL queries performed:

- Displaying the names of the unique launches sites in the space mission;
- Display 5 records where launch sites begin with 'CCA';
- Display the total payload mass carried by boosters launched by NASA;
- Display the average of the amount paid by booster version F9 v1.1;
- List of the dates of the first successful landing outcome;
- List of the names of the boosters which have success in drone ships and have payloads between 4000 and 6000;
- The total od number of successful and failure mission outcomes;
- List of the names of the *booster\_version* which have carried the maximum payload mass;
- Listing the records which will display the successful *landing\_outcomes* in the ground pad, booster versions, and launch site for the months in the year 2015;
- Ranking the count of successful *landing\_outcomes* between the date 2010-06-04 and 2017-03-20.



# Build an Interactive Map with Folium

- To visualize the Launch Data into an interactive map:
  - We took the Latitude and Longitude Coordinates at each launch site and added a Circle Marker around each launch site with a label with the name of the launch site
  - We assigned the dataframe launch\_outcomes and converted to classes 0 and 1 with **Green** and **Red**, respectively;
  - The red circles at each launch site coordinates with label showing launch site name;
- The objects are created to understand better the data problem. Also, we can easily show all launch sites and the successful and unsuccessful landings.





# Build a Dashboard with Plotly Dash

---

- The dashboard has a dropdown, pie chart and scatters plot:
  - The dropdown allows you to choose the launch site;
  - The pie chart shows the total of successful and unsuccessful launches sites selected from the dropdowns;
- Scatter Graph establishing the relationship with Outcome and Payload Mass (Kg) for the different Booster Versions:
  - It shows the relationship between two variables.
  - It is best to show you a non-linear pattern.
  - The range of data flow, maximum and minimum value, can be determined.
  - Observation and reading are straightforward





# Predictive Analysis (Classification)

---

## Data Preparation

- Load our dataset into NumPy and Pandas
- Data transformation
- Split our data into training and test data sets

## Model Preparation

- Check how many test samples we have
- Decide which type of machine learning algorithms we want to use
- Set our parameters and algorithms to GridSearchCV
- Training GridSearchCV

## Model Evaluation

- Check the accuracy of each model
- Get tuned hyperparameters for each type of algorithm

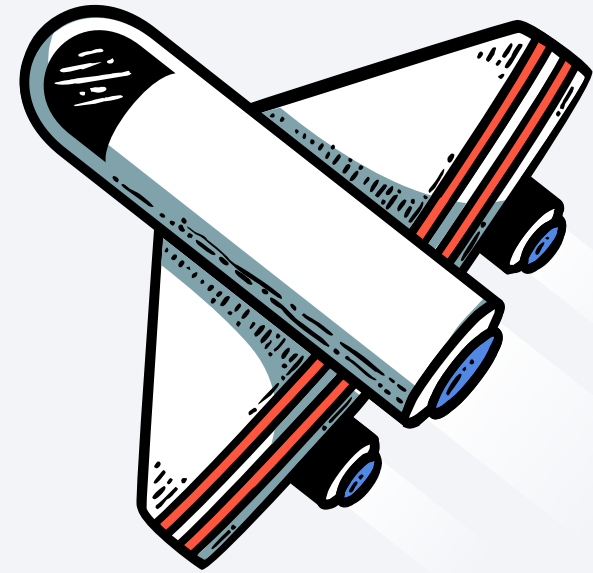
## Improving the model and finding the best classification model

- Feature Engineering
- Algorithm Tuning
- Comparison between methods
- The model with the best accuracy score wins the best-performing model

# Results

---

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





The background of the slide is an abstract composition. It features a dark blue base color. Overlaid on this are numerous diagonal streaks in shades of red and cyan. A faint, light blue grid pattern is also visible, particularly in the lower half of the image. The overall effect is dynamic and technological.

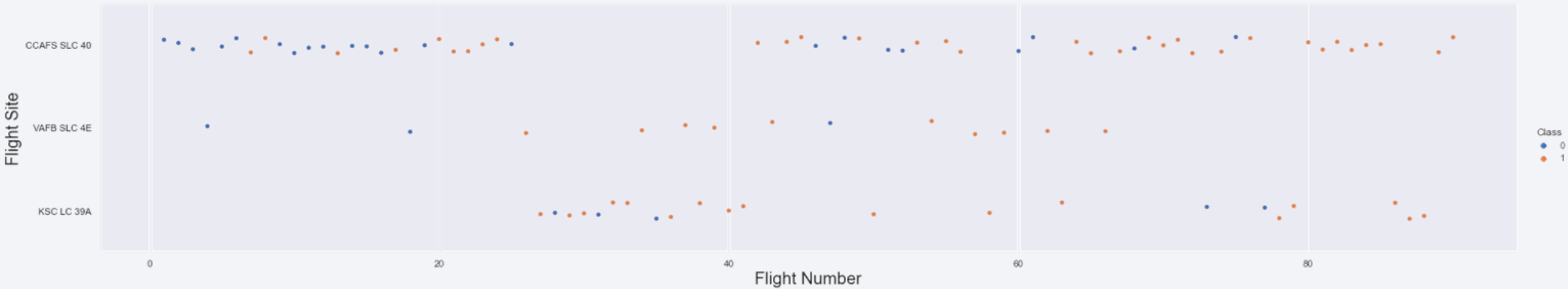
Section 2

# Insights drawn from EDA





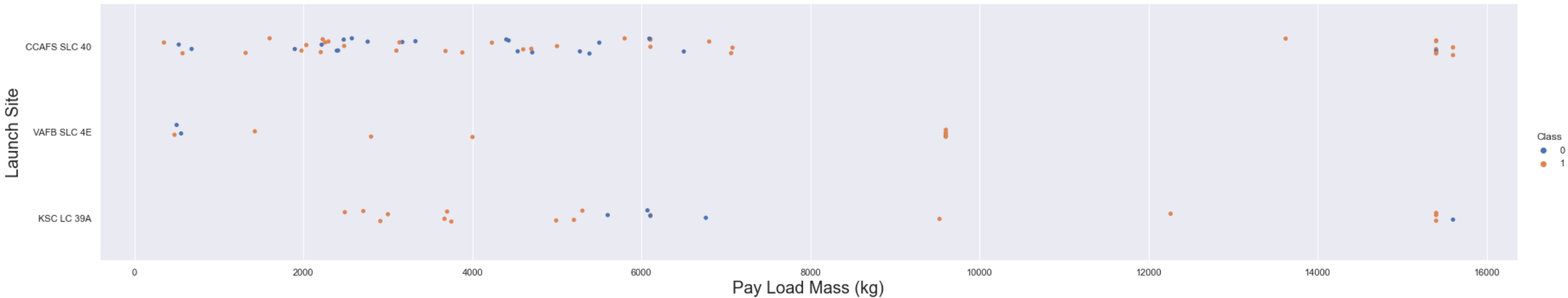
# Flight Number vs. Launch Site



It was observed, for each site, the success rate is increasing



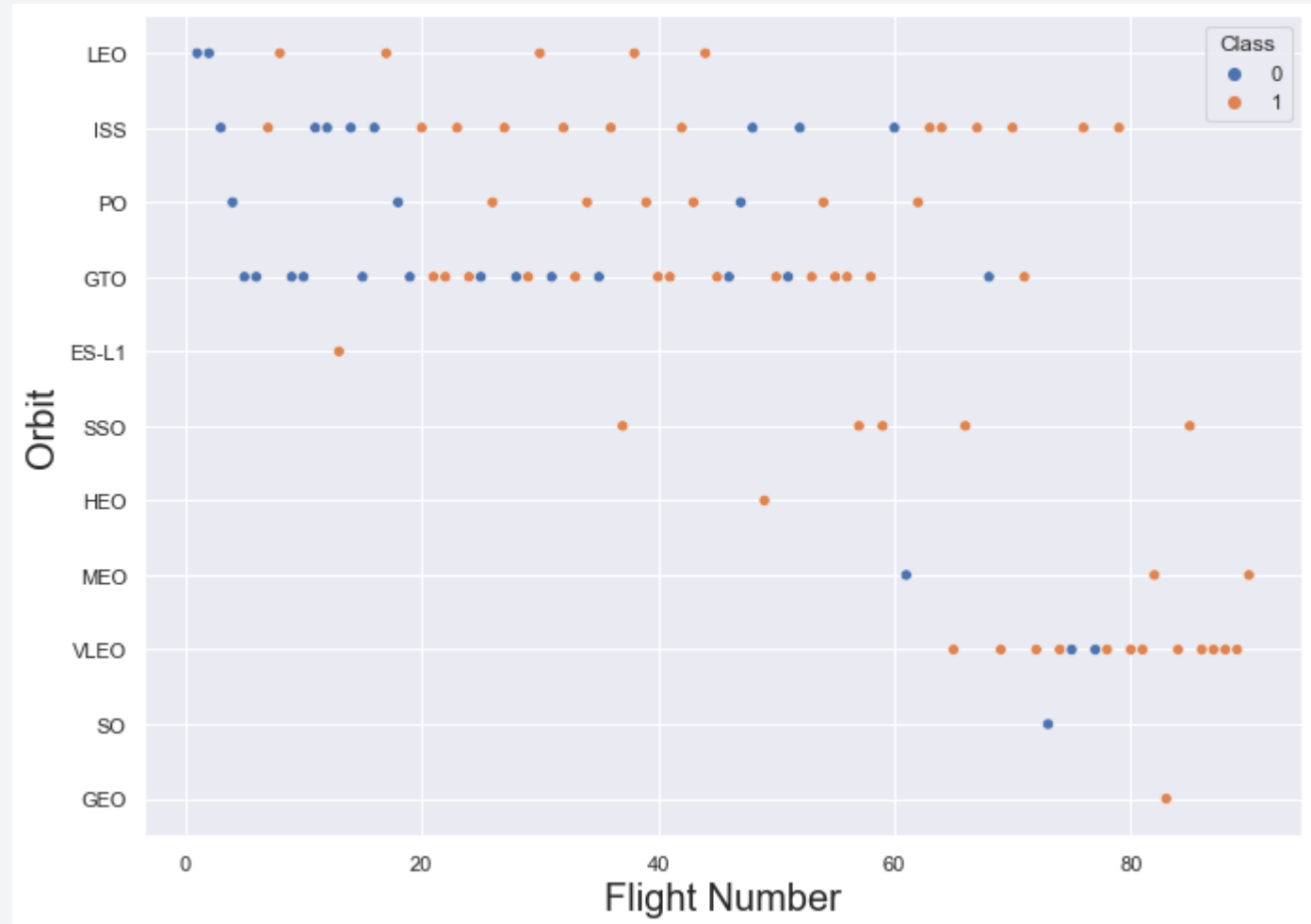
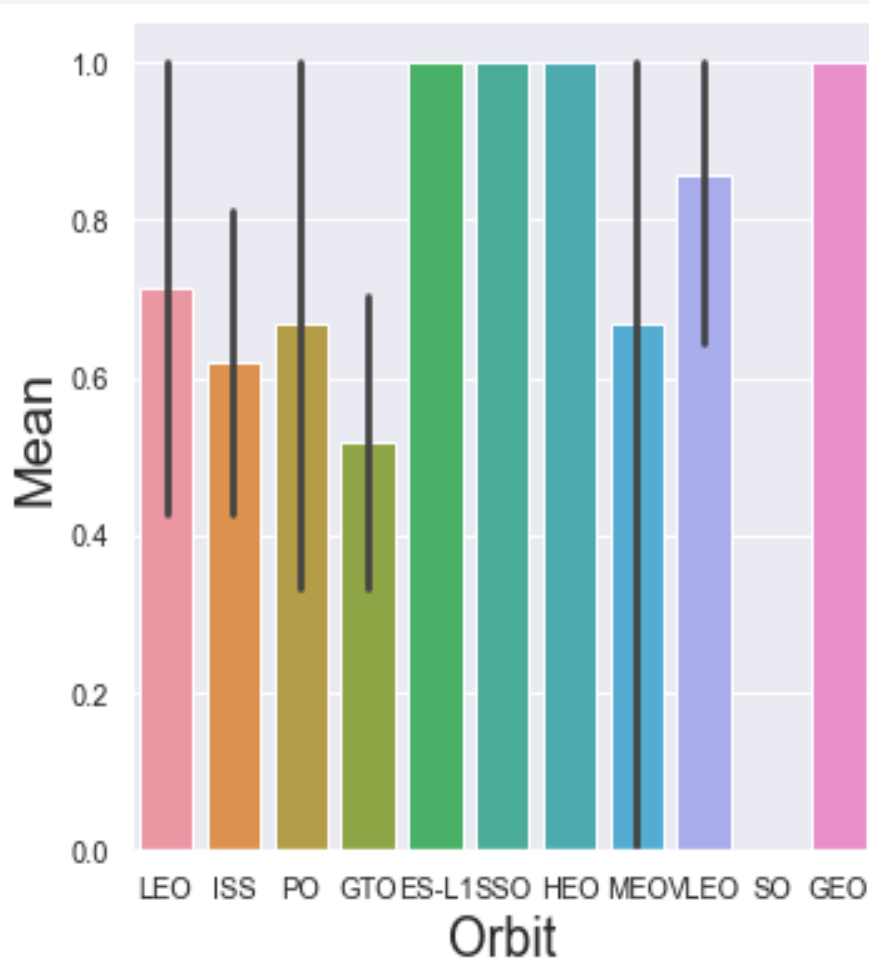
# Payload vs. Launch Site



The heave impacts the launch site. Therefore, a heavier payload may be a consideration for a successful landing.



# Success Rate vs. Orbit Type



We found that the success rate increases with the number of flights for the LEO orbit. On the other hand, in other orbits, like GTO, it's not related to success.



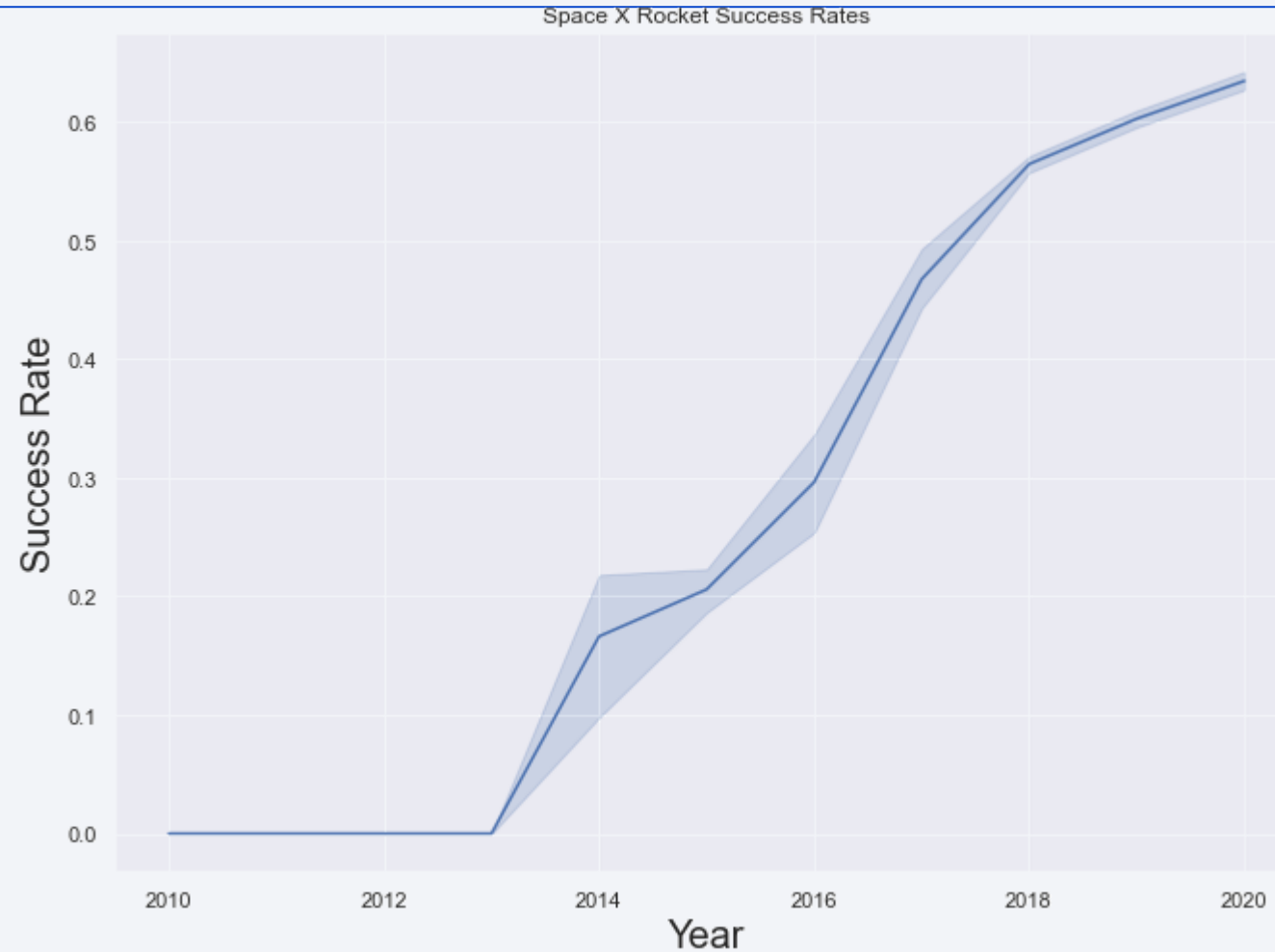
# Payload vs. Orbit Type



Heavy payloads have a negative influence on GTO orbits and positive on GTO and Polar LEO orbits.



# Launch Success Yearly Trend



Since 2013, the success rate has been increasing.





# All Launch Site Names

---

```
1 %sql select DISTINCT launch_site from SpaceX
```

Python

launch\_site

CCAFS LC-40

CCAFS SLC-40

KSC LC-39A

VAFB SLC-4E

To use DISTINCT in a query remove all duplicated values



# Launch Site Names Begin with 'CCA'

```
1 %sql select * from SpaceX WHERE launch_site LIKE 'CCA%' limit 5
```

Python

DATE	time_utc_	booster_version	launch_site	payload	payload_mass_kg_	orbit	customer
2010-06-04	18:45:00	F9 v1.0 B0003	CCAFS LC-40	Dragon Spacecraft Qualification Unit	0	LEO	SpaceX
2010-12-08	15:43:00	F9 v1.0 B0004	CCAFS LC-40	Dragon demo flight C1, two CubeSats, barrel of Brouere cheese	0	LEO (ISS)	NASA (COTS) NRO
2012-05-22	07:44:00	F9 v1.0 B0005	CCAFS LC-40	Dragon demo flight C2	525	LEO (ISS)	NASA (COTS)
2012-10-08	00:35:00	F9 v1.0 B0006	CCAFS LC-40	SpaceX CRS-1	500	LEO (ISS)	NASA (CRS)
2013-03-01	15:10:00	F9 v1.0 B0007	CCAFS LC-40	SpaceX CRS-2	677	LEO (ISS)	NASA (CRS)

To use WHERE followed by LIKE allows get the launches that contain subsisting 'CCA', and getting only the 5 rows using LIMIT 5



# Total Payload Mass

---

```
1 %sql select SUM(payload_mass__kg_) from SpaceX where Customer = 'NASA (CRS)'
```

Python

```
</> 1
45596
```

The query gives the sum of all payload where the customer is equal to NASA (CRS)



# Average Payload Mass by F9 v1.1

---

```
1 %sql select AVG(payload_mass__kg_) from SpaceX where Booster_Version = 'F9 v1.1'
```

```
</> 1  
2928
```

The WHERE clause filters the dataset to only perform calculations on Booster\_version F9 v1.1



# First Successful Ground Landing Date

```
1 %sql select MIN(DATE) from SpaceX where landing_outcome = 'Success (ground pad)'
```

Python

</>

1

2015-12-22

The WHERE clause filters the dataset to only perform calculations on Landing\_Outcome = Success (ground pad) and get the first day



# Successful Drone Ship Landing with Payload between 4000 and 6000

---

```
1 %sql select DISTINCT booster_version from SpaceX where landing__outcome = 'Success'
```

Python

```
</> booster_version  
F9 FT B1021.2  
F9 FT B1031.2  
F9 FT B1022  
F9 FT B1026
```

[Click here - Notebook Link](#)



The query returns the booster version where landing was successful and payload was between 4k and 6k.

# Total Number of Successful and Failure Mission Outcomes

---

```
1 %sql SELECT mission_outcome, COUNT(*) FROM SpaceX GROUP BY mission_outcome
```

	mission_outcome	2
	Failure (in flight)	1
	Success	99
	Success (payload status unclear)	1

[Click here - Notebook Link](#)



Using the function COUNT works out the amount

# Boosters Carried Maximum Payload

---

```
</> booster_version  
F9 B5 B1048.4  
F9 B5 B1049.4  
F9 B5 B1051.3  
F9 B5 B1056.4  
F9 B5 B1048.5  
F9 B5 B1051.4  
F9 B5 B1049.5  
F9 B5 B1060.2  
F9 B5 B1058.3  
F9 B5 B1051.6  
F9 B5 B1060.3  
F9 B5 B1049.7
```

Using the word **SELECT** in the query means that it will show values in the **Booster\_Version** column from Space

[Click here - Notebook Link](#)





# 2015 Launch Records

---

DATE	landing_outcome	booster_version	launch_site
2015-01-10	Failure (drone ship)	F9 v1.1 B1012	CCAFS LC-40
2015-04-14	Failure (drone ship)	F9 v1.1 B1015	CCAFS LC-40

DATE LIKE puts the value of 2015

# Rank Landing Outcomes Between 2010-06-04 and 2017-03-20

</>	DATE	time_utc_	booster_version	launch_site	payload	payload_mass_kg_	orbit	customer	mission_outcome	landing_outcome
	2016-06-15	14:29:00	F9 FT B1024	CCAFS LC-40	ABS-2A Eutelsat 117 West B	3600	GTO	ABS Eutelsat	Success	Failure (drone ship)
	2016-03-04	23:35:00	F9 FT B1020	CCAFS LC-40	SES-9	5271	GTO	SES	Success	Failure (drone ship)
	2016-01-17	18:42:00	F9 v1.1 B1017	VAFB SLC-4E	Jason-3	553	LEO	NASA (LSP) NOAA CNES	Success	Failure (drone ship)
	2015-04-14	20:10:00	F9 v1.1 B1015	CCAFS LC-40	SpaceX CRS-6	1898	LEO (ISS)	NASA (CRS)	Success	Failure (drone ship)
	2015-01-10	09:47:00	F9 v1.1 B1012	CCAFS LC-40	SpaceX CRS-5	2395	LEO (ISS)	NASA (CRS)	Success	Failure (drone ship)
	2010-12-08	15:43:00	F9 v1.0 B0004	CCAFS LC-40	Dragon demo flight C1, two CubeSats, barrel of Brouere cheese	0	LEO (ISS)	NASA (COTS) NRO	Success	Failure (parachute)
	2010-06-04	18:45:00	F9 v1.0 B0003	CCAFS LC-40	Dragon Spacecraft Qualification Unit	0	LEO	SpaceX	Success	Failure (parachute)

Function WHERE filters landing\_\_outcome and LIKE (Success or Failure; AND (DATE between) DESC means its arranging the dataset into descending order

[Click here - Notebook Link](#)





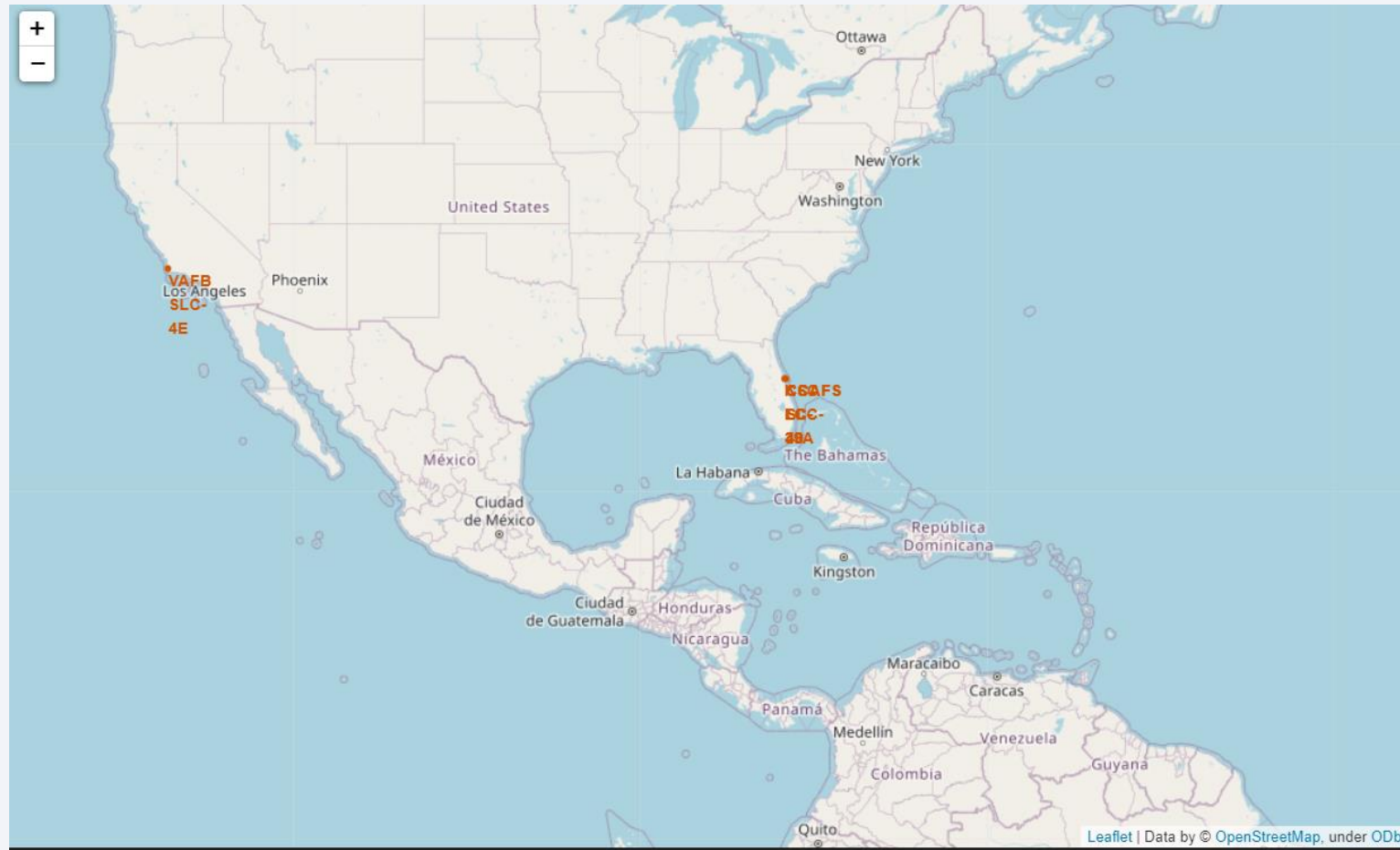
A satellite view of Earth from space, showing the curvature of the planet and city lights at night. The background is a deep blue gradient.

Section 3

# Launch Sites Proximities Analysis



# SpaceX in the US



SpaceX in the US



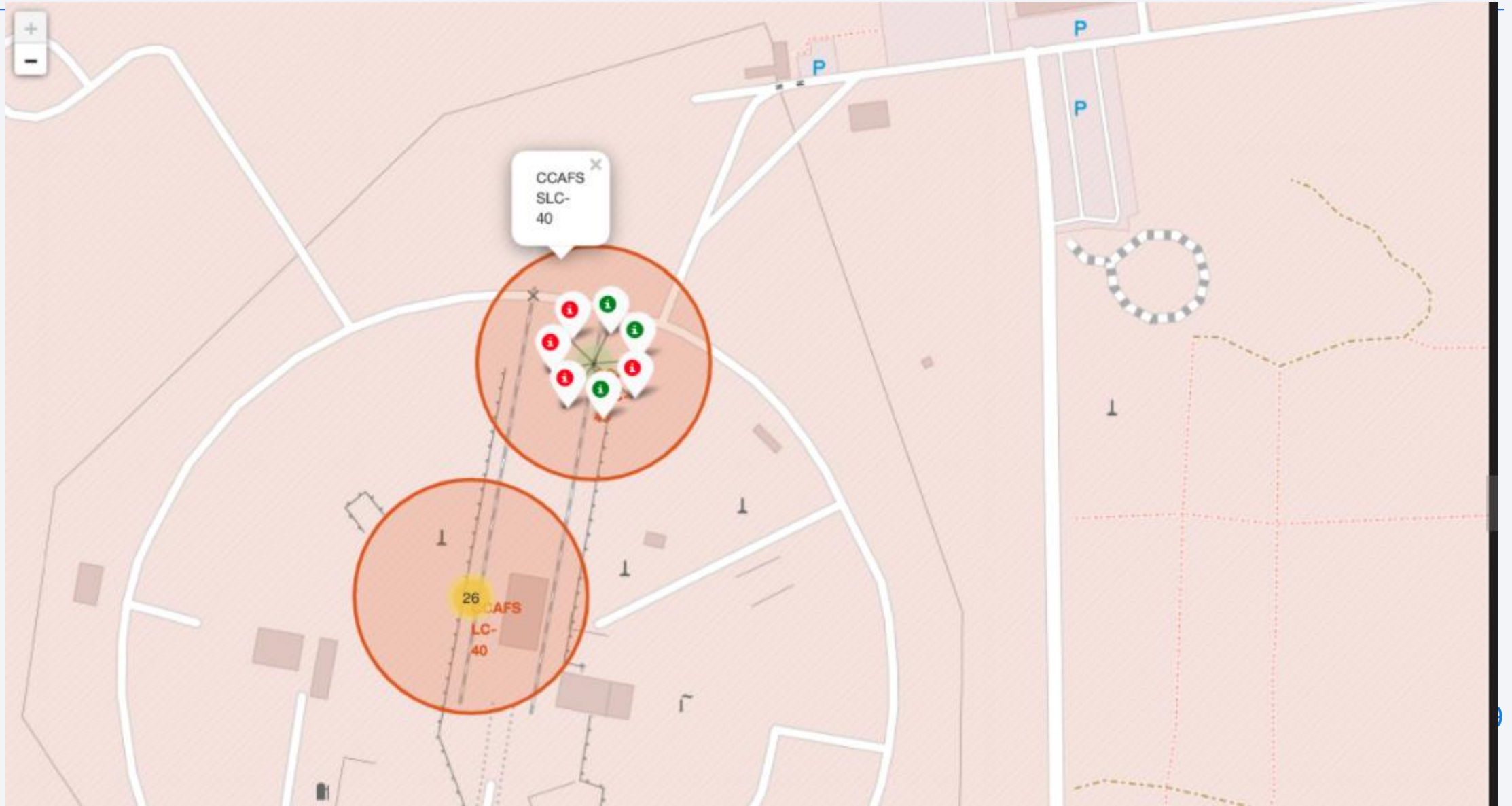
# Color-labeled launch outcomes







# Color-labeled launch





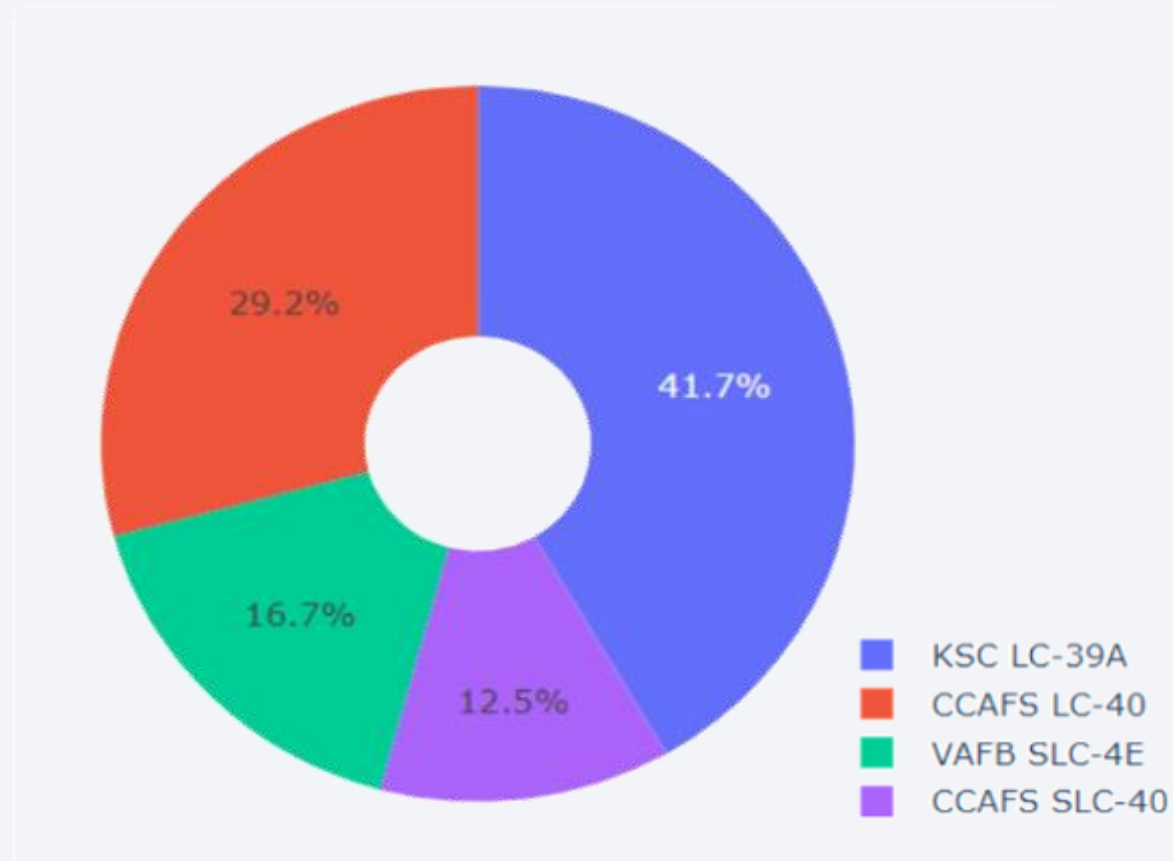


Section 4

# Build a Dashboard with Plotly Dash

# Pie chart showing the success percentage

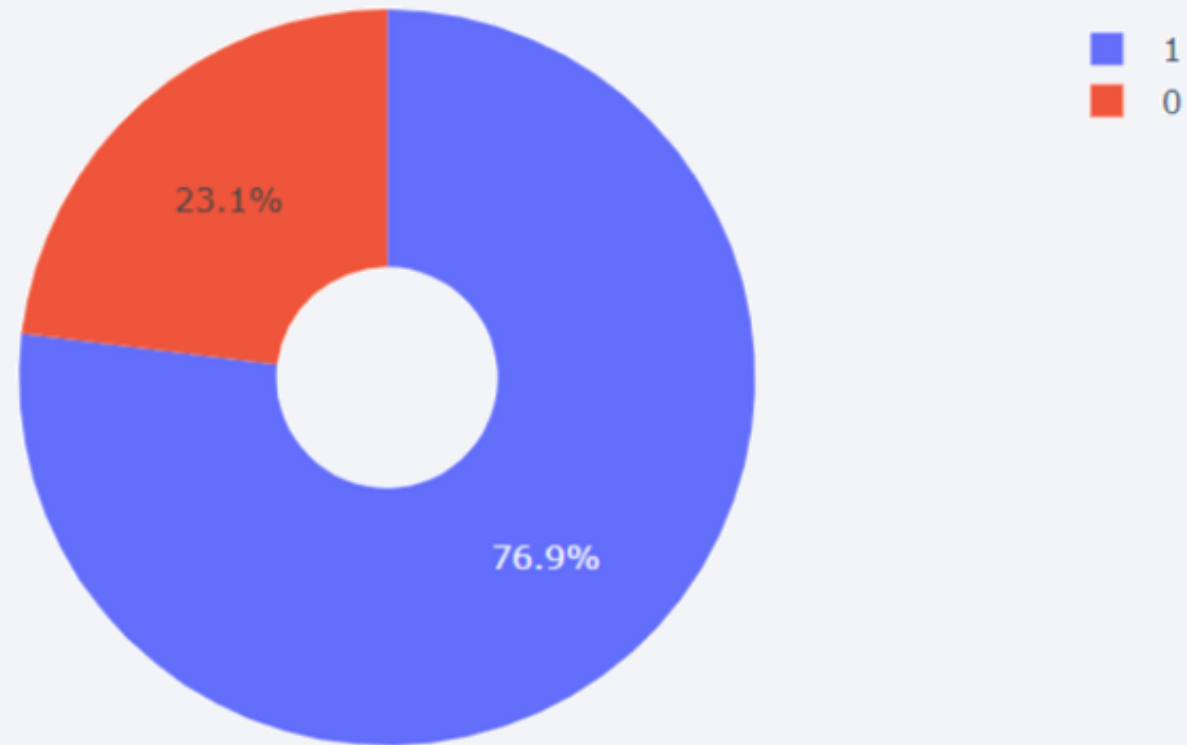
---



KSC had the most successful launches.

## Pie chart showing the Launch site with the highest launch success ratio

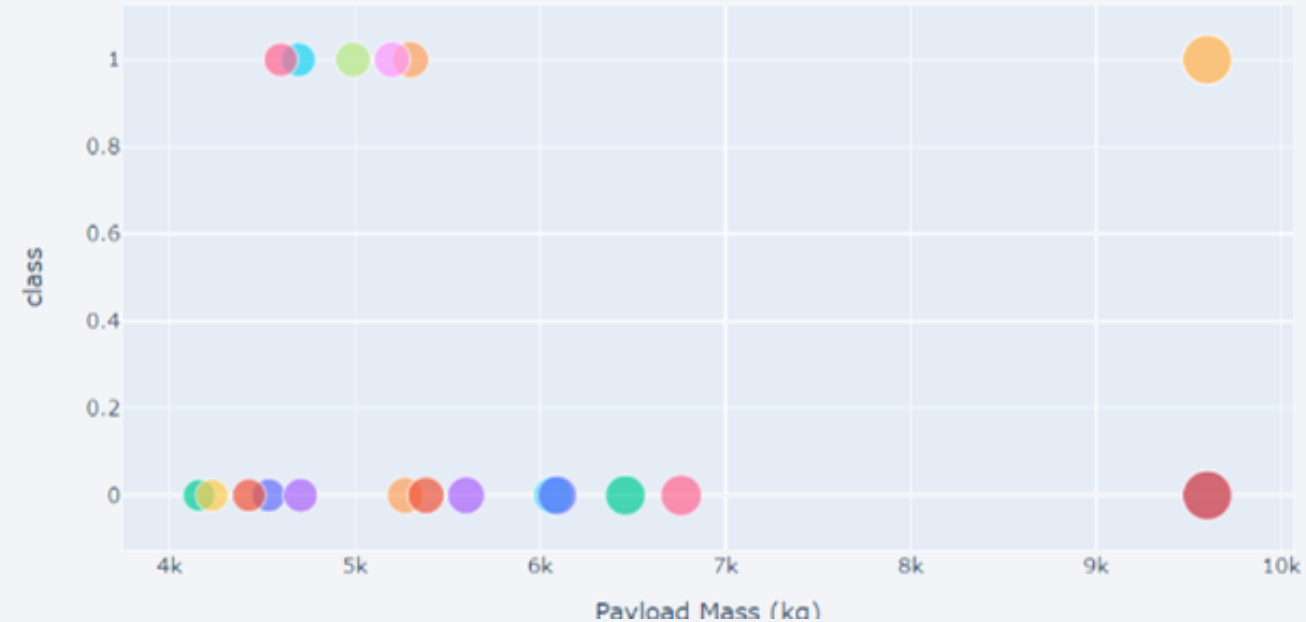
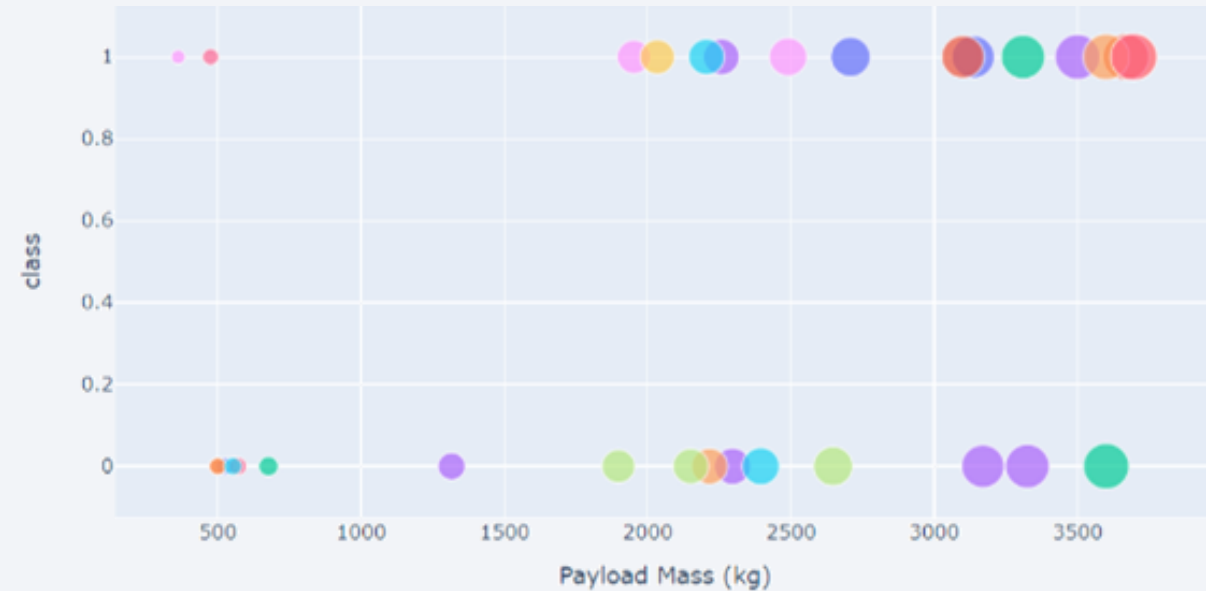
---



**KSC LC-39A achieved a 76,9% of success.**

## Scatter plot of Payload vs Launch Outcome for all sites, with different payload selected in the range slider

---



**The success rates for low weighted payload is higher than the heavy weighted payloads.**

Section 5

# Predictive Analysis (Classification)





# Classification Accuracy

```
1 parameters = {'criterion': ['gini', 'entropy'],
2               'splitter': ['best', 'random'],
3               'max_depth': [2*n for n in range(1,10)],
4               'max_features': ['auto', 'sqrt'],
5               'min_samples_leaf': [1, 2, 4],
6               'min_samples_split': [2, 5, 10]}
7
8 tree = DecisionTreeClassifier()

[31] ✓ 0.1s Python

1 tree_cv=GridSearchCV(tree, param_grid=parameters, cv=10)
2 tree_cv.fit(X_train,Y_train)

[32] ✓ 6.2s Python

... GridSearchCV(cv=10, estimator=DecisionTreeClassifier(),
               param_grid={'criterion': ['gini', 'entropy'],
                           'max_depth': [2, 4, 6, 8, 10, 12, 14, 16, 18],
                           'max_features': ['auto', 'sqrt'],
                           'min_samples_leaf': [1, 2, 4],
                           'min_samples_split': [2, 5, 10],
                           'splitter': ['best', 'random']})

1 print("tuned hyperparameters :(best parameters) ",tree_cv.best_params_)
2 print("accuracy :",tree_cv.best_score_)

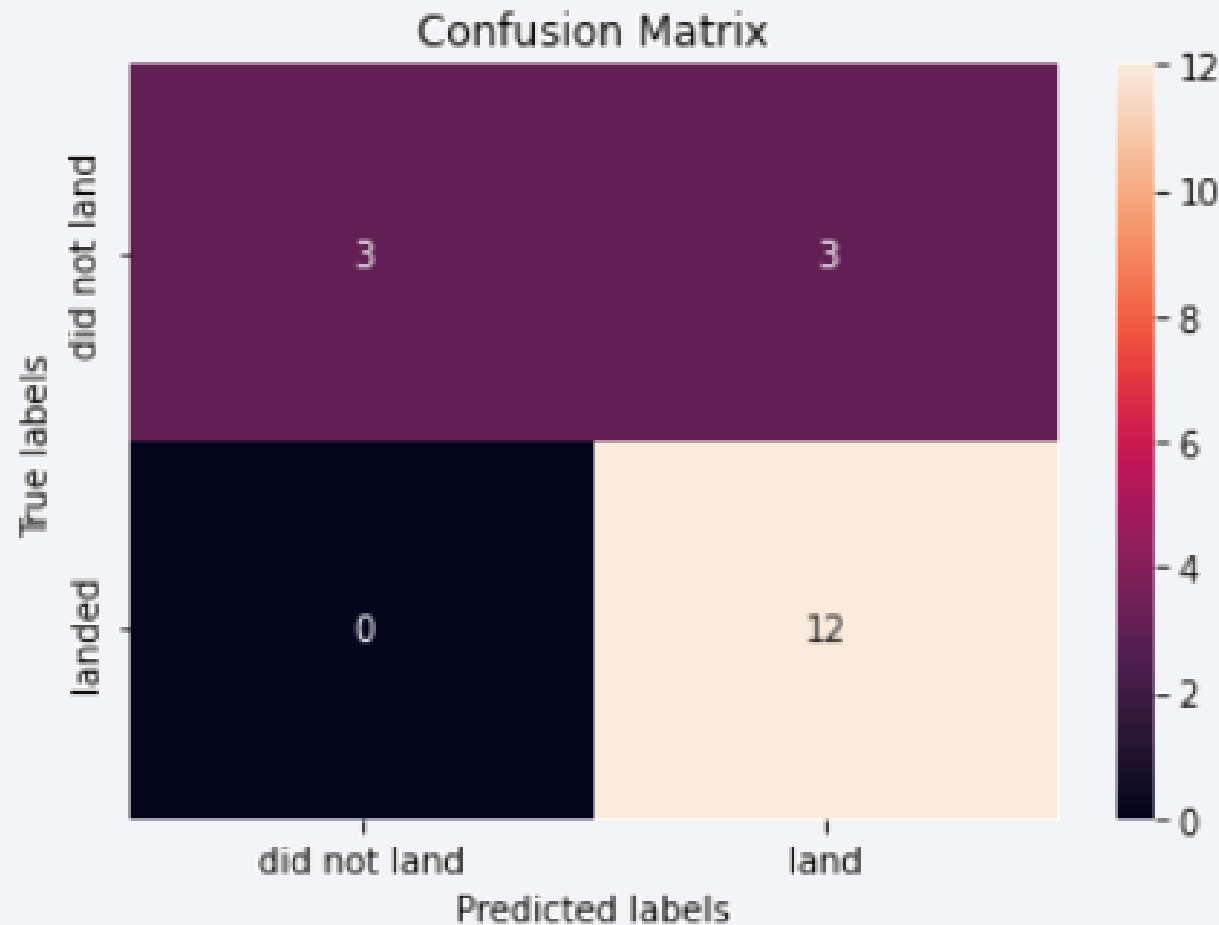
[33] ✓ 0.7s Python

... tuned hyperparameters :(best parameters) {'criterion': 'gini', 'max_depth': 4, 'max_features': 'sqrt', 'min_samples_leaf': 1,
'min_samples_split': 2, 'splitter': 'random'}
accuracy : 0.8767857142857143
```

The decision tree was the best model based in the classification accuracy.



# Confusion Matrix



The confusion matrix for the decision tree classifier shows that the classifier can distinguish between the classes.



# Conclusions

---

- In conclusion, we can assume:
  - The larger the flight amount at a launch site;
  - The success rates for SpaceX launches is directly proportional time in years;
  - Orbits ES-L1, GEO, HEO, SSO, VLEO had the most success rate
  - The Decision tree classifier is the best machine learning algorithm for this task;

Thank you!

