



Outline

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

Executive summary

Methods Summary

- Data collection and Data wrangling
 - Exploratory Data Analysis (EDA) with Data Visualization
 - EDA with SQL
- Building an interactive map with Folium
- Dashboard with Plotly Dash
- Predictive analysis

Results summary

- Exploratory Data Analysis results
- Interactive maps and dashboard
- Predictive results

Introduction

■ Context and Background of Project

- **Project aim:** Predict the successful landing of the Falcon 9 first stage.
- **SpaceX's claim:** Falcon 9 rocket launch costs \$62 million, significantly lower than competitors.
- **Cost comparison:** Other providers charge upwards of \$165 million for similar launches.
- **SpaceX's advantage:** Ability to reuse the first stage, reducing costs.
- **Importance of prediction:** Determines the cost of a launch based on successful landing.
- **Relevance to competitors:** Valuable information for companies aiming to compete with SpaceX in rocket launches.

■ Problems need to focus on

- What factors influences the successful landing of rocket?
- How do specific relationships with rocket variables affect the success rate of landing?
- What conditions does SpaceX need to achieve to optimize results and ensure the highest success rate for rocket landings?

Methodology

- **Executive Summary**
- Data collection methodology:
 - SpaceX REST API
 - Web Scrapping from Wikipedia
- Perform data wrangling
 - Dropping unnecessary columns
 - One Hot Encoding for classification models
- Perform exploratory data analysis (EDA) using visualization and SQL
- Perform interactive visual analytics using Folium and Plotly Dash
- Perform predictive analysis using classification models

Methodology

- The dataset was obtained by accessing SpaceX launch data through the SpaceX REST API.
- This API offers detailed information on launches, including rocket details, payload, launch specifications, landing specifications, and landing outcomes.
- The primary goal is to utilize this dataset to predict whether SpaceX will attempt to land a rocket.
- The SpaceX REST API endpoints begin with `api.spacexdata.com/v4/`.
- An alternative method for acquiring Falcon 9 Launch data involves web scraping Wikipedia using BeautifulSoup.

Data Collection – SpaceX API

Getting Response from API

```
spacex_url="https://api.spacexdata.com/v4/launches/past"  
response = requests.get(spacex_url)
```

Convert Response to JSON File

```
data = response.json()  
data = pd.json_normalize(data)
```

Create dataframe

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

Create dictionary with data

```
launch_dict = {'FlightNumber': list(data['flight_number']),  
'Date': list(data['date']),  
'BoosterVersion': BoosterVersion,  
'PayloadMass': PayloadMass,  
'Orbit': Orbit,  
'LaunchSite': LaunchSite,  
'Outcome': Outcome,  
'Flights': Flights,  
'GridFins': GridFins,  
'Reused': Reused,  
'Legs': Legs,  
'LandingPad': LandingPad,  
'Block': Block,  
'ReusedCount': ReusedCount,  
'Serial': Serial,  
'Longitude': Longitude,  
'Latitude': Latitude}
```

Transform data

```
getLaunchSite(data)  
getPayloadData(data)  
getCoreData(data)  
getBoosterVersion(data)
```

Filter dataframe

```
data_falcon9 = data[data['BoosterVersion']!='Falcon 1']
```

Export to file

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

Data Collection – Web Scrapping

Getting Response from HTML

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

Create BeautifulSoup Object

```
soup = BeautifulSoup(response.text, "html5lib")
```

Create dictionary

```
launch_dict= dict.fromkeys(column_names)

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

# Let's initial the launch_dict with each value to be an empty list
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'] = []
# Added some new columns
launch_dict['Version Booster']=[]
launch_dict['Booster landing']=[]
launch_dict['Date']=[]
launch_dict['Time']=[]
```

Get column names

```
for th in first_launch_table.find_all('th'):
    name = extract_column_from_header(th)
    if name is not None and len(name) > 0 :
        column_names.append(name)
```

Find all tables

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

Add data to keys

```
extracted_row = 0
#Extract each table
for table_number,table in enumerate(soup.find_all
    # get table row
    for rows in table.find_all("tr"):
        #check to see if first table heading is a
        if rows.th:
            if rows.th.string:
                flight_number=rows.th.string.stri
                flag=flight_number.isdigit()
```

Create dataframe from dictionary

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

Export to file

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


Data Wrangling

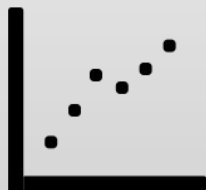
- The dataset contains various instances where booster landings were not successful, categorized based on different criteria.
- Instances where a landing was attempted but failed due to an accident are identified.
- Two main categories denote the landing outcomes: "Ocean" and "RTLS" (Return to Launch Site).
- Under the "Ocean" category, "True Ocean" indicates a successful landing in a specific ocean region, while "False Ocean" indicates an unsuccessful landing in the ocean.
- In the "RTLS" category, "True RTLS" signifies a successful landing on a ground pad at the launch site, while "False RTLS" indicates an unsuccessful landing on a ground pad.
- Additionally, the dataset includes outcomes categorized as "ASDS" (Autonomous Spaceport Drone Ship).
- Within the "ASDS" category, "True ASDS" represents a successful landing on a drone ship, whereas "False ASDS" denotes an unsuccessful landing on a drone ship.
- These outcomes are transformed into training labels for analysis, where the successful landing is labeled as '1' and the unsuccessful landing as '0'.

EDA with Data Visualization

- **Scatter Graphs**

- 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

- *Scatter plots show relationship between*
- *variables. This relationship is called the correlation.*



- **Bar Graph**

- Success rate vs. Orbit

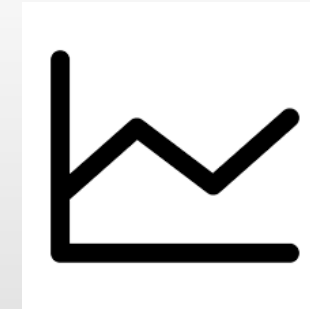
Bar graphs show the relationship between numeric and categorical variables.



- **Line Graph**

- Success rate vs. Year

Line graphs show data variables and their trends. Line graphs can help to show global behavior and make prediction for unseen data.



EDA with SQL

- We performed SQL queries to gather and understand data from dataset:
 - Displaying the names of the unique launch sites in the space mission.
 - Display 5 records where launch sites begin with the string 'CCA'
 - Display the total payload mass carried by boosters launched by NASA (CRS).
 - Display average payload mass carried by booster version F9 v1.1.
 - List the date when the first successful landing outcome in ground pad was achieved.
 - List the names of the boosters which have success in drone ship and have payload mass greater than 4000 but less than 6000.
 - List the total number of successful and failure mission outcomes.
 - List the names of the booster versions which have carried the maximum payload mass.
 - List the records which will display the month names, failure landing_outcomes in drone ship, booster versions, launch_site for the months in year 2015.
 - Rank the count of successful landing outcomes between the date 04-06-2010 and 20-03-2017 in descending order.

Interactive map with Folium

- Folium map object is a map centered on NASA Johnson Space Center at Houston, Texas
 - Red circle at NASA Johnson Space Center's coordinate with label showing its name (*folium.Circle, folium.map.Marker*).
 - Red circles at each launch site coordinates with label showing launch site name (*folium.Circle, folium.map.Marker, folium.features.DivIcon*).
 - The grouping of points in a cluster to display multiple and different information for the same coordinates (*folium.plugins.MarkerCluster*).
 - Markers to show successful and unsuccessful landings. Green for successful landing and Red for unsuccessful landing. (*folium.map.Marker, folium.Icon*).
 - Markers to show distance between launch site to key locations (railway, highway, coastway, city) and plot a line between them. (*folium.map.Marker, folium.PolyLine, folium.features.DivIcon*)
- These objects are created in order to understand better the problem and the data. We can show easily all launch sites, their surroundings and the number of successful and unsuccessful landings.

Dashboard with Plotly Dash

- Dashboard has dropdown, pie chart, rangeslider and scatter plot components
 - Dropdown allows a user to choose the launch site or all launch sites (*dash_core_components.Dropdown*).
 - Pie chart shows the total success and the total failure for the launch site chosen with the dropdown component (*plotly.express.pie*).
 - Rangeslider allows a user to select a payload mass in a fixed range (*dash_core_components.RangeSlider*).
 - Scatter chart shows the relationship between two variables, in particular Success vs Payload Mass (*plotly.express.scatter*).

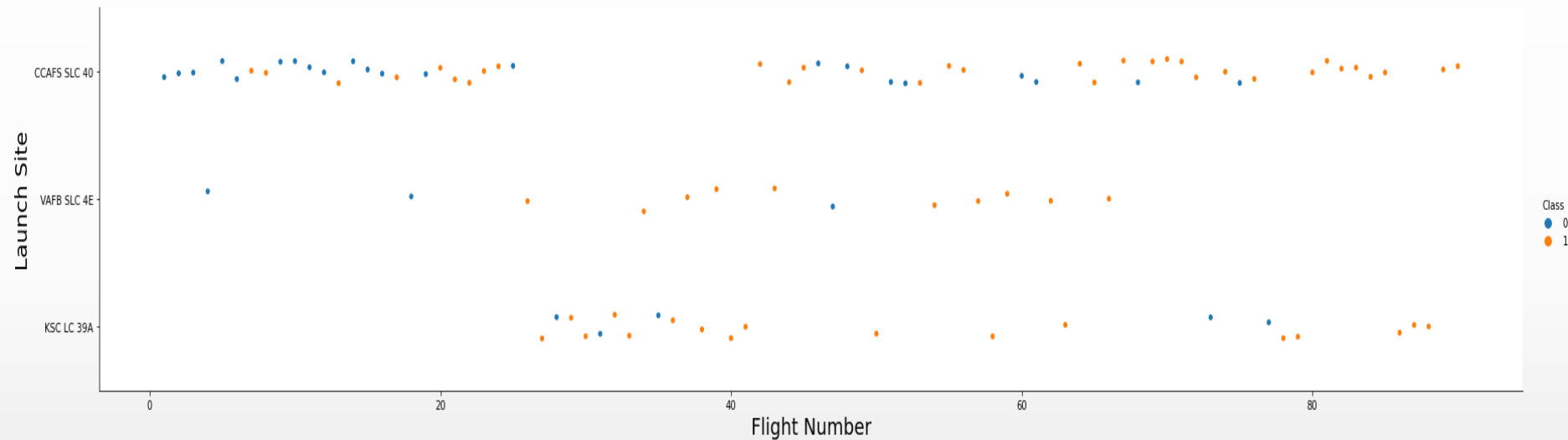
Predictive Analysis

- Data preparation
 - Load dataset
 - Normalize data
 - Split data into training and test sets.
- Model preparation
 - Selection of machine learning algorithms
 - Set parameters for each algorithm to GridSearchCV
 - Training GridSearchModel models with training dataset
- Model evaluation
 - Get best hyperparameters for each type of model
 - Compute accuracy for each model with test dataset
 - Plot Confusion Matrix
- Model comparison
 - Comparison of models according to their accuracy
 - The model with the best accuracy will be chosen (see Notebook for result)

Results

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

Flight Number vs. Launch site



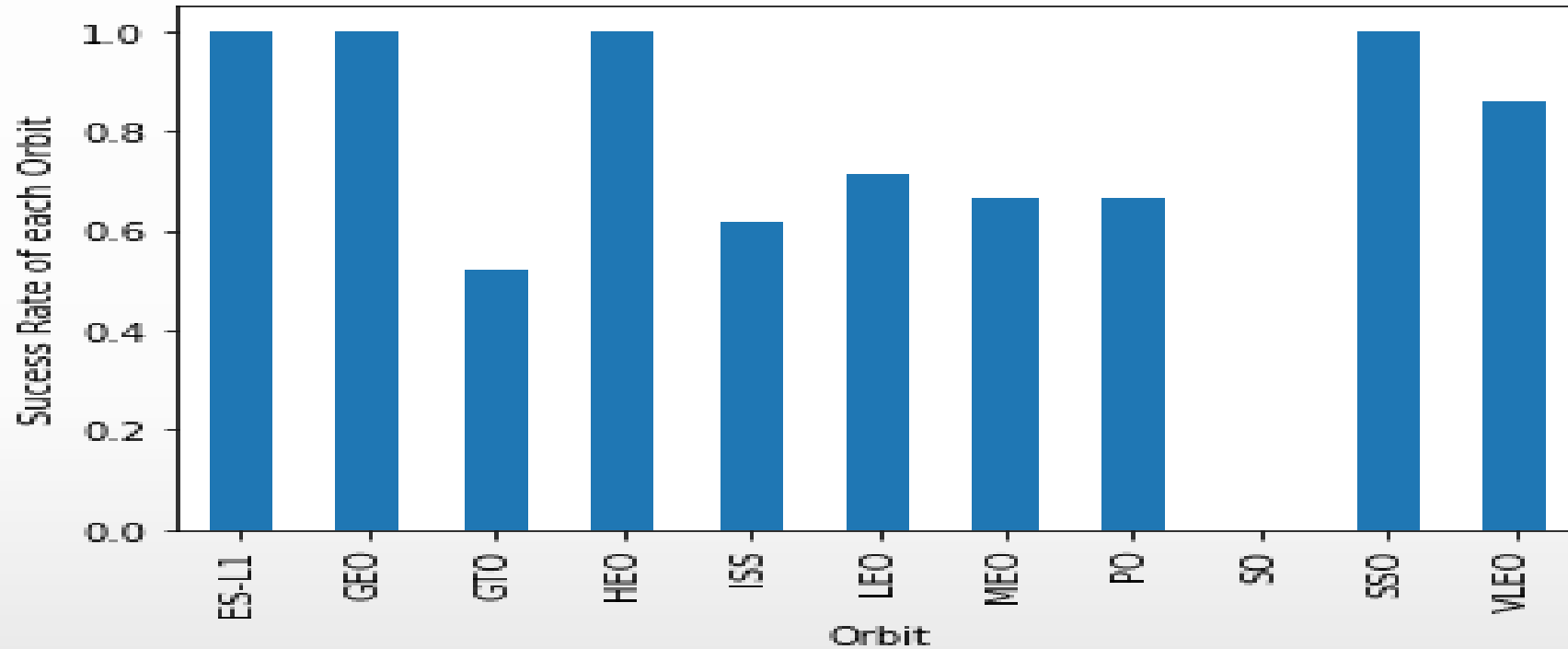
Observation: for each site, the success rate is increasing.

Payload vs. Launch site



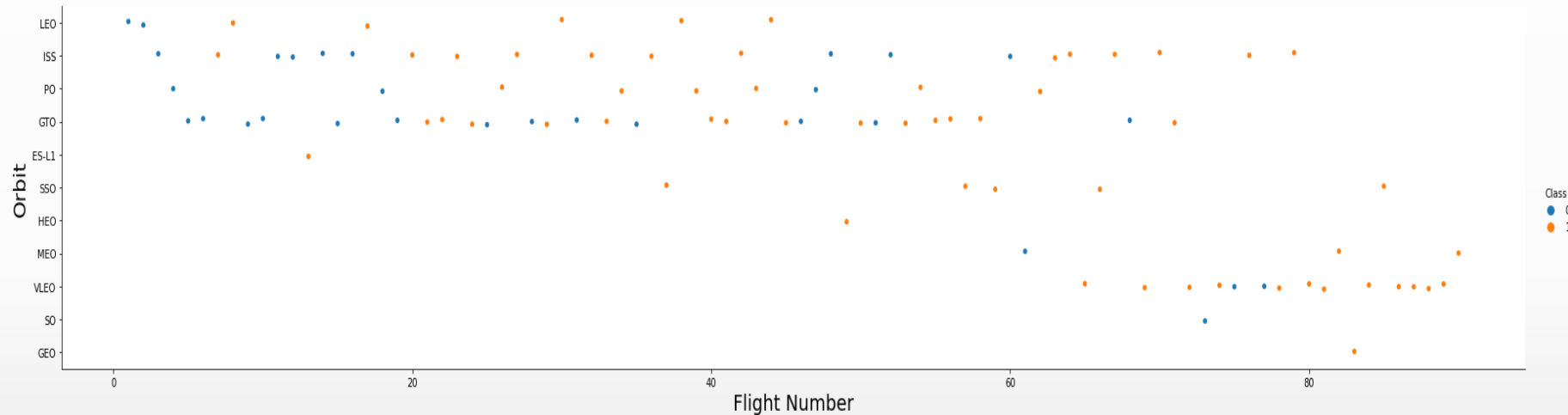
Depending on the launch site, a heavier payload may be a consideration for a successful landing. On the other hand, a too heavy payload can make a landing fail.

Success vs. Orbit type



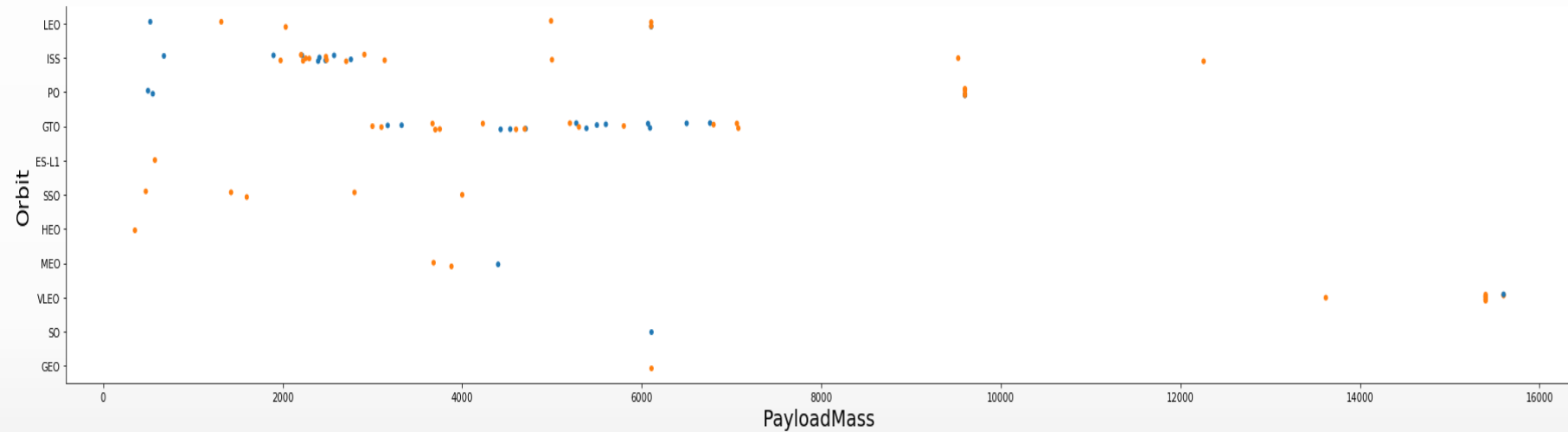
With this plot, we can see success rate for different orbit types. We note that ES-L1, GEO, HEO, SSO have the best success rate.

Flight number vs. orbit type



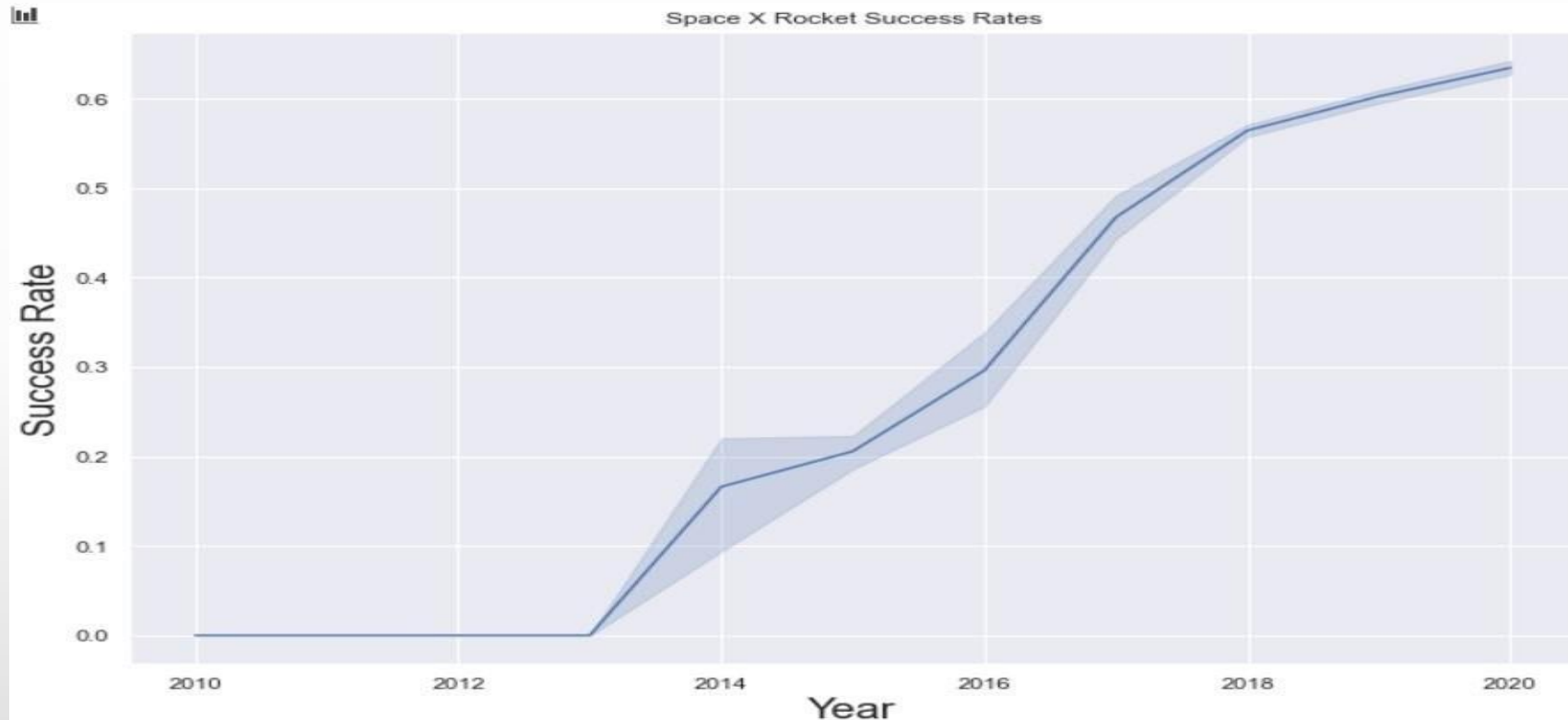
Observation: The success rate increases with the number of flights for the LEO orbit. For some orbits like GTO, there is no relation between the success rate and the number of flights. But we can suppose that the high success rate of some orbits like SSO or HEO is due to the knowledge learned during former launches for other orbits.

Payload vs. orbit type



The weight of the payloads can have a great influence on the success rate of the launches in certain orbits. For example, heavier payloads improve the success rate for the LEO orbit. Another finding is that decreasing the payload weight for a GTO orbit improves the success of a launch.

Launch success yearly trend



Observation: Since 2013, we can see an increase in the Space X Rocket success rate.

All Launch sites names

SQLQuery

```
SELECT DISTINCT "LAUNCH_SITE" FROM SPACEXTBL
```

Explanation

The use of DISTINCT in the query allows to remove duplicate LAUNCH_SITE.

Results

Launch_Site
CCAFS LC-40
VAFB SLC-4E
KSC LC-39A
CCAFS SLC-40

Launch site names begin with CCA

SQLQuery

```
SELECT * FROM SPACEXTBL WHERE "LAUNCH_SITE" LIKE '%CA%' LIMIT 5
```

Results

Explanation

The WHERE clause followed by LIKE clause filters launch sites that contain the substring CCA. LIMIT 5 shows 5 records from filtering.

Date	Time (UTC)	Booster_Version	Launch_Site	Payload	PAYLOAD_MASS_KG_	Orbit	Customer
04-06-2010	18:45:00	F9 v1.0 B0003	CCAFS LC-40	Dragon Spacecraft Qualification Unit	0	LEO	SpaceX
08-12-2010	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
22-05-2012	07:44:00	F9 v1.0 B0005	CCAFS LC-40	Dragon demo flight C2	525	LEO (ISS)	NASA (COTS)
08-10-2012	00:35:00	F9 v1.0 B0006	CCAFS LC-40	SpaceX CRS-1	500	LEO (ISS)	NASA (CRS)
01-03-2013	15:10:00	F9 v1.0 B0007	CCAFS LC-40	SpaceX CRS-2	677	LEO (ISS)	NASA (CRS)

Total Payload mass

SQLQuery

```
SELECT SUM("PAYLOAD_MASS_KG_") FROM SPACEXTBL WHERE "CUSTOMER" = 'NASA (CRS)'
```

Explanation

This query returns the sum of all payload masses where the customer is NASA (CRS).

Results

SUM("PAYLOAD_MASS_KG_")
45596

Average Payload Mass by F9 v1.1

SQLQuery

```
SELECT AVG("PAYLOAD_MASS_KG_") FROM SPACEXTBL WHERE "BOOSTER_VERSION" LIKE '%F9 v1.1%'
```

Explanation

This query returns the average of all payload masses where the booster version contains the substring F9 v1.1.

Results

AVG("PAYLOAD_MASS_KG_")

2534.6666666666665

First Successful Ground Landing Date

SQLQuery

```
SELECT MIN("DATE") FROM SPACEXTBL WHERE "Landing _Outcome" LIKE '%Success%'
```

Explanation

With this query, we select the oldest successful landing.

The WHERE clause filters dataset in order to keep only records where landing was successful. With the MIN function, we select the record with the oldest date.

Results

```
MIN("DATE")
```

```
01-05-2017
```

Successful Drone Ship Landing with Payload between 4000 and 6000

SQLQuery

```
%sql SELECT "BOOSTER_VERSION" FROM SPACEXTBL WHERE "LANDING_OUTCOME" = 'Success (drone ship)' \
AND "PAYLOAD_MASS_KG_" > 4000 AND "PAYLOAD_MASS_KG_" < 6000;
```

Explanation

This query returns the booster version where landing was successful and payload mass is between 4000 and 6000 kg. The WHERE and AND clauses filter the dataset.

Results

Booster_Version
F9 FT B1022
F9 FT B1026
F9 FT B1021.2
F9 FT B1031.2

Total Number of Successful and Failure Mission Outcomes

SQLQuery

```
%sql SELECT (SELECT COUNT("MISSION_OUTCOME") FROM SPACEXTBL WHERE "MISSION_OUTCOME" LIKE '%Success%') AS SUCCESS, \
(SELECT COUNT("MISSION_OUTCOME") FROM SPACEXTBL WHERE "MISSION_OUTCOME" LIKE '%Failure%') AS FAILURE
```

Explanation

With the first SELECT, we show the subqueries that return results. The first subquery counts the successful mission. The second subquery counts the unsuccessful mission. The WHERE clause followed by LIKE clause filters mission outcome. The COUNT function counts records filtered.

Results

SUCCESS	FAILURE
100	1

Boosters Carried Maximum Payload

SQL Query

```
%sql SELECT DISTINCT "BOOSTER_VERSION" FROM SPACEXTBL \
WHERE "PAYLOAD_MASS_KG_" = (SELECT max("PAYLOAD_MASS_KG_") FROM SPACEXTBL)
```

Explanation

We used a subquery to filter data by returning only the heaviest payload mass with MAX function. The main query uses subquery results and returns unique booster version (SELECT DISTINCT) with the heaviest payload mass.

Results

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

2015 Launch Records

SQL Query

```
%sql SELECT substr("DATE", 4, 2) AS MONTH, "BOOSTER_VERSION", "LAUNCH_SITE" FROM SPACEXTBL\
WHERE "LANDING_OUTCOME" = 'Failure (drone ship)' and substr("DATE",7,4) = '2015'
```

Explanation

This query returns month, booster version, launch site where landing was unsuccessful and landing date took place in 2015. Substr function process date in order to take month or year. Substr(DATE, 4, 2) shows month. Substr(DATE, 7, 4) shows year.

Results

MONTH	Booster_Version	Launch_Site
01	F9 v1.1 B1012	CCAFS LC-40
04	F9 v1.1 B1015	CCAFS LC-40

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

SQLQuery

```
%sql SELECT "LANDING_OUTCOME", COUNT("LANDING_OUTCOME") FROM SPACEXTBL\
WHERE "DATE" >= '04-06-2010' and "DATE" <= '20-03-2017' and "LANDING_OUTCOME" LIKE '%Success%\
GROUP BY "LANDING_OUTCOME" \
ORDER BY COUNT("LANDING_OUTCOME") DESC ;
```

Explanation

This query returns landing outcomes and their count where mission was successful and date is between 04/06/2010 and 20/03/2017. The group by clause groups results by landing outcome and order by countedesc shows results in decreasing order.

Results

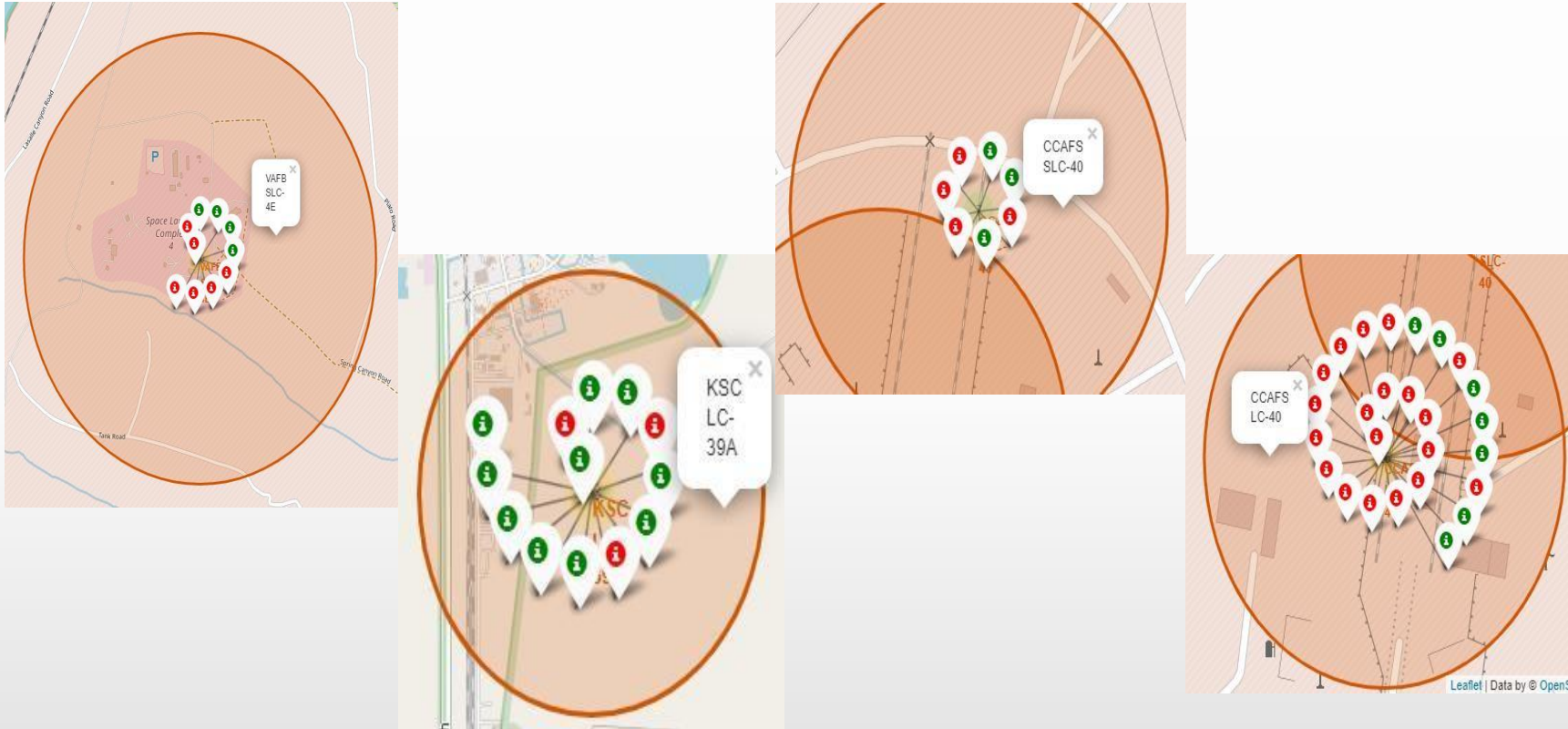
Landing_Outcome	COUNT("LANDING_OUTCOME")
Success	20
Success (drone ship)	8
Success (ground pad)	6

All launch site global map marker

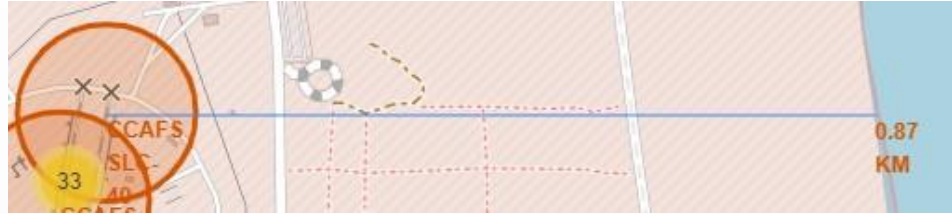


We see that Space X launch sites are located on the coast of the United States

Color labelled Markers



Folium Map – Distances between CCAFS SLC-40 and its proximities



Is CCAFS SLC-40 in close proximity to railways ? Yes
Is CCAFS SLC-40 in close proximity to highways ? Yes
Is CCAFS SLC-40 in close proximity to coastline ?
Yes

Do CCAFS SLC-40 keeps certain distance away from cities ? No

Dashboard – Total success by Site

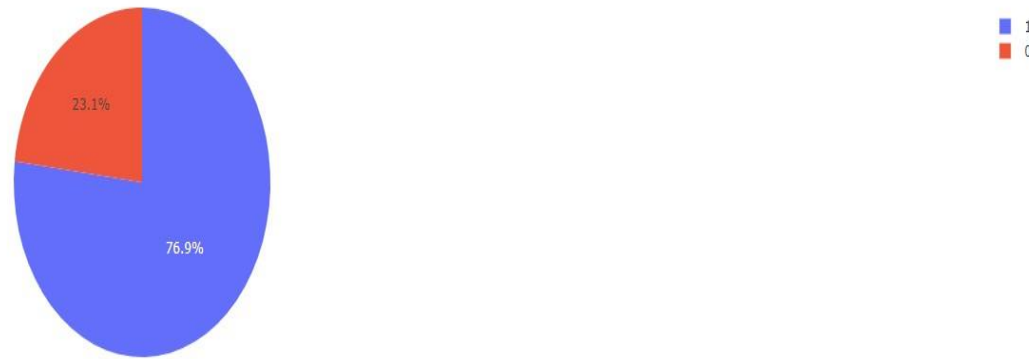
Total Success Launches by Site



We see that KSC LC-39A has the best success rate of launches.

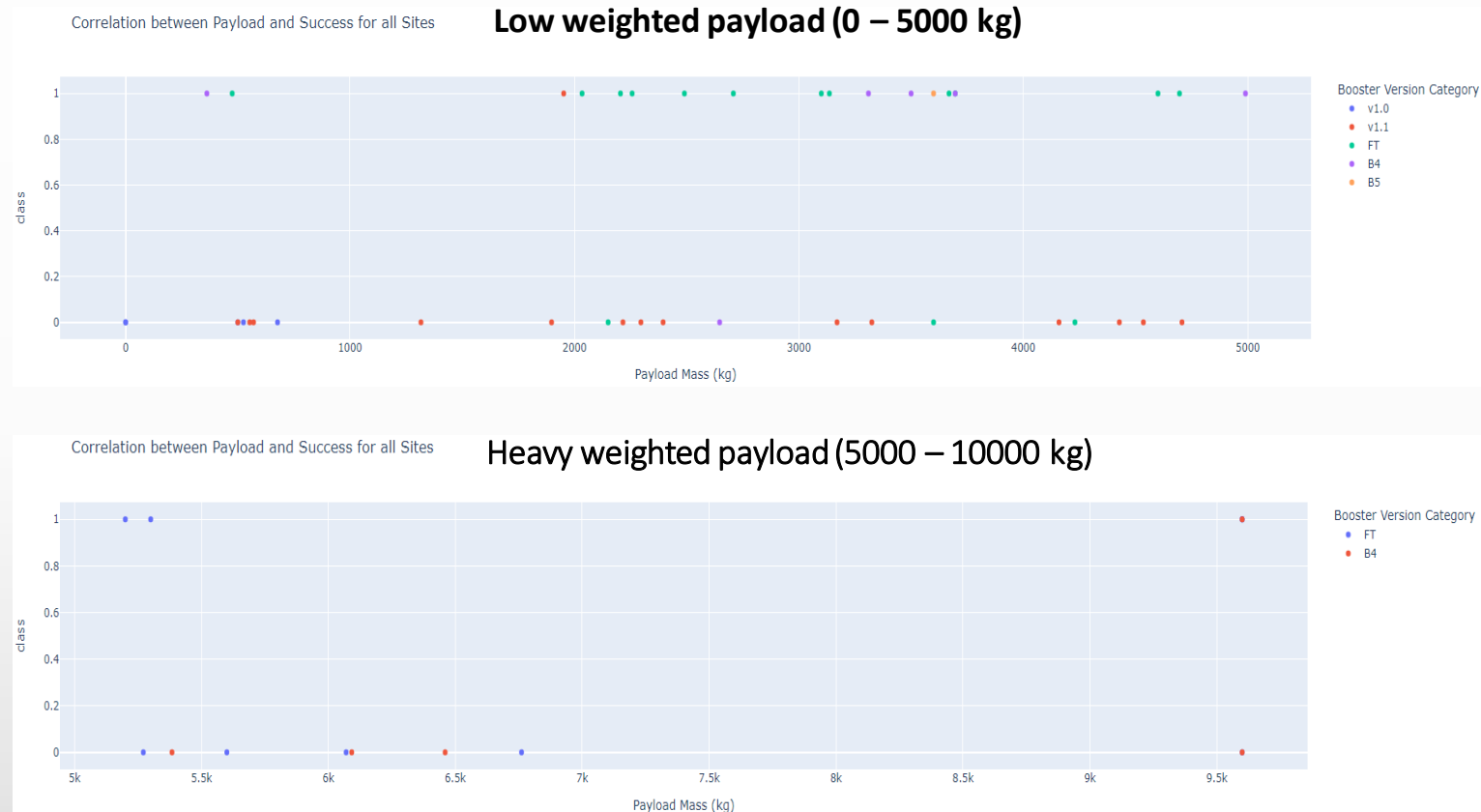
Dashboard – Total success launches for Site KSC LC-39A

Total Success Launches for Site KSC LC-39A



We see that KSC LC-39A has achieved a 76.9% success rate while getting a 23.1% failure rate.

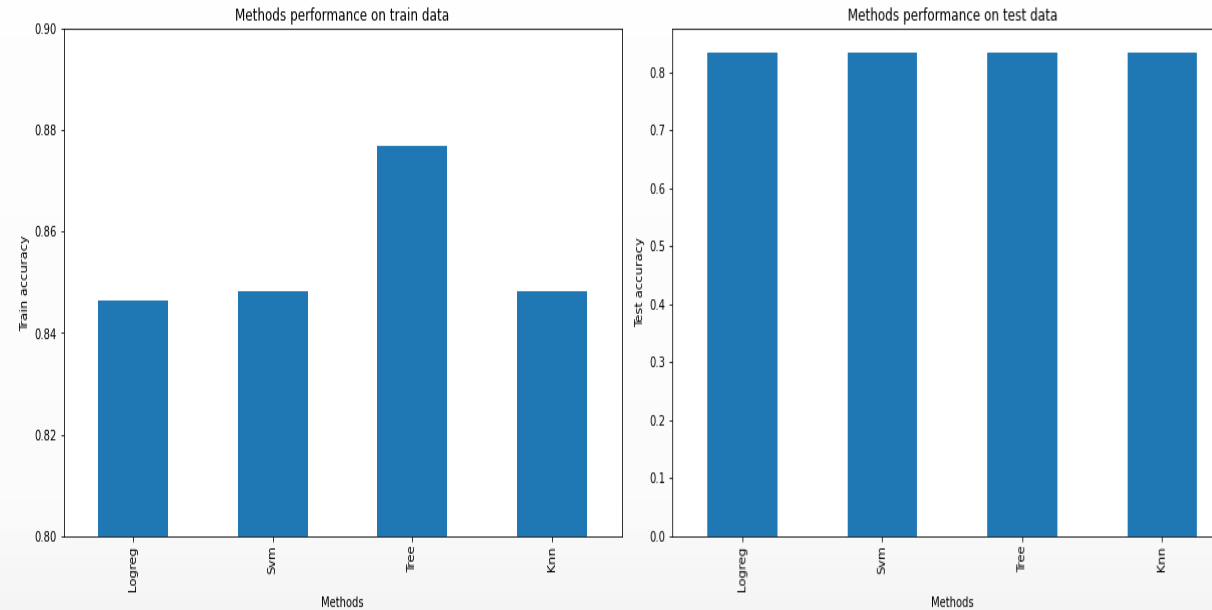
Dashboard — Payload mass vs Outcome for all sites with different payload mass selected



Low weighted payloads have a better success rate than the heavy weighted payloads.

Classification Accuracy

	Accuracy Train	Accuracy Test
Tree	0.876786	0.833333
Knn	0.848214	0.833333
Svm	0.848214	0.833333
Logreg	0.846429	0.833333



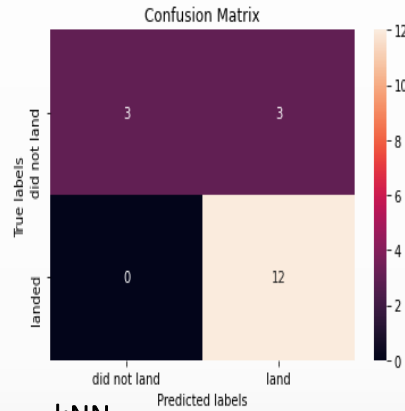
For accuracy test, all methods performed similar. We could get more test data to decide between them. But if we really need to choose one right now, we would take the decision tree.

Decision tree best parameters

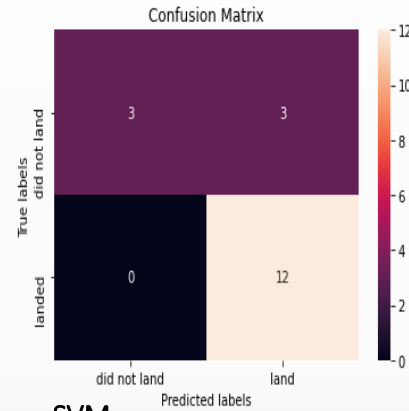
```
tuned hyperparameters :(best parameters) {'criterion': 'entropy', 'max_depth': 12, 'max_features': 'sqrt', 'min_samples_leaf': 4, 'min_samples_split': 2, 'splitter': 'random'}
```

Confusion Matrix

Logistic regression

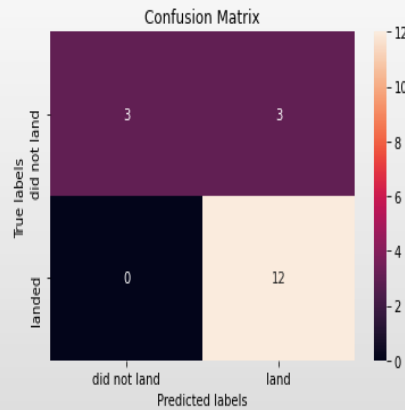


Decision Tree

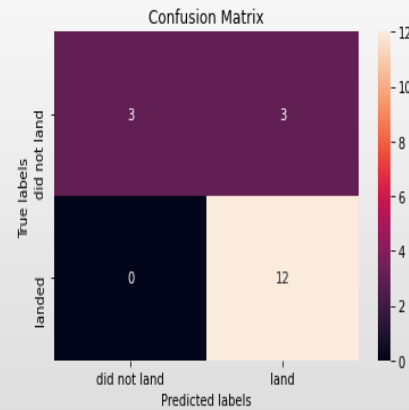


As the test accuracy are all equal, the confusion matrices are also identical. The main problem of these models are false positives.

kNN



SVM



		Actual values	
		1	0
Predicted values	1	TP	FP
	0	FN	TN

Conclusions

- The success of a mission can be explained by several factors such as the launch site, the orbit and especially the number of previous launches. Indeed, we can assume that there has been a gain in knowledge between launches that allowed to go from a launch failure to a success.
- The orbits with the best success rates are GEO, HEO, SSO, ES-L1.
- Depending on the orbits, the payload mass can be a criterion to take into account for the success of a mission. Some orbits require a light or heavy payload mass. But generally low weighted payloads perform better than the heavy weighted payloads.
- With the current data, we cannot explain why some launch sites are better than others (KSCLC-39A is the best launch site). To get an answer to this problem, we could obtain atmospheric or other relevant data.
- For this dataset, we choose the Decision Tree Algorithm as the best model even if the test accuracy between all the models used is identical. We choose Decision Tree Algorithm because it has a better train accuracy.