

Applied Data Science Project- Capstone

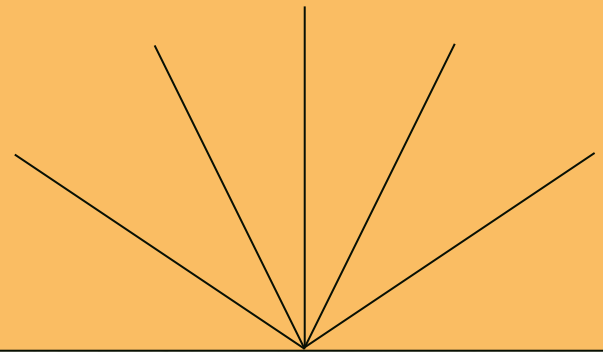
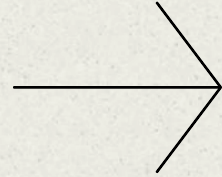


TABLE OF CONTENTS

○1 Executive
Summary

○3 Methodology

○2 Introduction

○4 Conclusion

Executive Summary

01.



Summary of methodologies

Data Collection via API, Web
Scraping

Exploratory Data Analysis (EDA)
with Data Visualization

EDA with SQL

Interactive Map with Folium

Dashboards with Plotly Dash

Predictive Analysis



Summary of all results

Exploratory Data Analysis results

Interactive maps and dashboard

Predictive results

O2.

Introduction

A solid orange horizontal bar spanning the width of the slide at the bottom.

Project background and context

- What are the main characteristics of a successful or failed landing ?
- What are the effects of each relationship of the rocket variables on the success or failure of a landing ?
- What are the conditions which will allow SpaceX to achieve the best landing success rate ?

Problems you want to find answers

The project aims to predict if the Falcon 9 first stage will successfully land. SpaceX charges \$62 million per launch, compared to other providers' \$165 million, due to their reusable first stage. Predicting landing success helps estimate launch costs, which is valuable for companies competing with SpaceX.

- 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

O3. Metho- dology

- Datasets are collected from Rest SpaceX API and webscrapping Wikipedia
 - The information obtained by the API are rocket, launches, payload information.
 - The Space X REST API URL is api.spacexdata.com/v4/
- The information obtained by the webscrapping of Wikipedia are launches, landing, payload information.
 - URL: https://en.wikipedia.org/w/index.php?title=List_of_Falcon_9_and_Falcon_Heavy_launches&oldid=1027686922

O4. Data Collection

Data Collection – SpaceX API

1. Getting Response from API

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

2. Convert Response to JSON File

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

3. Transform data

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

4. 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}
```

5. Create dataframe

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

6. Filter dataframe

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

7. Export to file

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

Data Collection – Scrapping

1. Getting Response from HTML

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

2. Create BeautifulSoup Object

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

3. Find all tables

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

4. 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)
```

5. 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'] = []
```

6. 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.strip()
                flag = flight_number.isdigit()
```

7. Create dataframe from dictionary

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

8. Export to file

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

Data Wrangling

- In the dataset, there are several cases where the booster did not land successfully.
 - True Ocean, True RTLS, True ASDS means the mission has been successful.
 - False Ocean, False RTLS, False ASDS means the mission was a failure.
- We need to transform string variables into categorical variables where 1 means the mission has been successful and 0 means the mission was a failure.

1. Calculate launches number for each site

```
df['LaunchSite'].value_counts()
CCAFS SLC 40    55
KSC LC 39A     22
VAFB SLC 4E     13
Name: LaunchSite, dtype: int64
```

2. Calculate the number and occurrence of each orbit

```
df['Orbit'].value_counts()
GTO    27
ISS    21
VLEO   14
PO      9
LEO      7
SSO      5
HEO      3
SO       1
ES-L1    1
HEO       1
GEO       1
Name: Orbit, dtype: int64
```

3. Calculate number

```
landing_outcomes = df['Outcome'].value_counts()
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
```

4. Create landing

```
landing_class = []
for key,value in df['Outcome'].items():
    if value in bad_outcomes:
        landing_class.append(0)
    else:
        landing_class.append(1)
df['Class']=landing_class
```

5. Export to file

```
df.to_csv("dataset_part_2.csv", index=False)
```

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 categoric 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.

Build an 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.

Build a 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*).

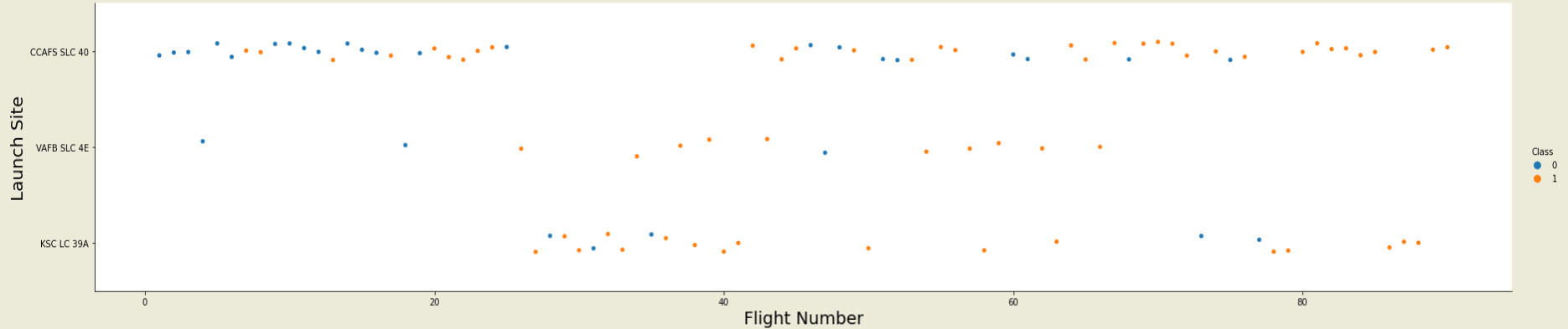
Predictive Analysis (Classification)

- 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

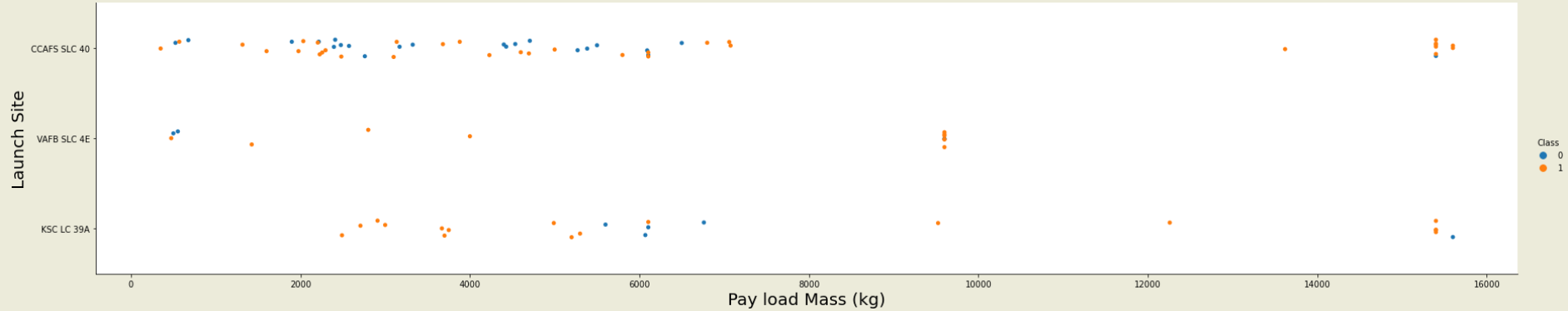
- Exploratory data analysis results
- Interactive analytics demo in screenshots
- Predictive analysis results

Flight Number & Launch Site



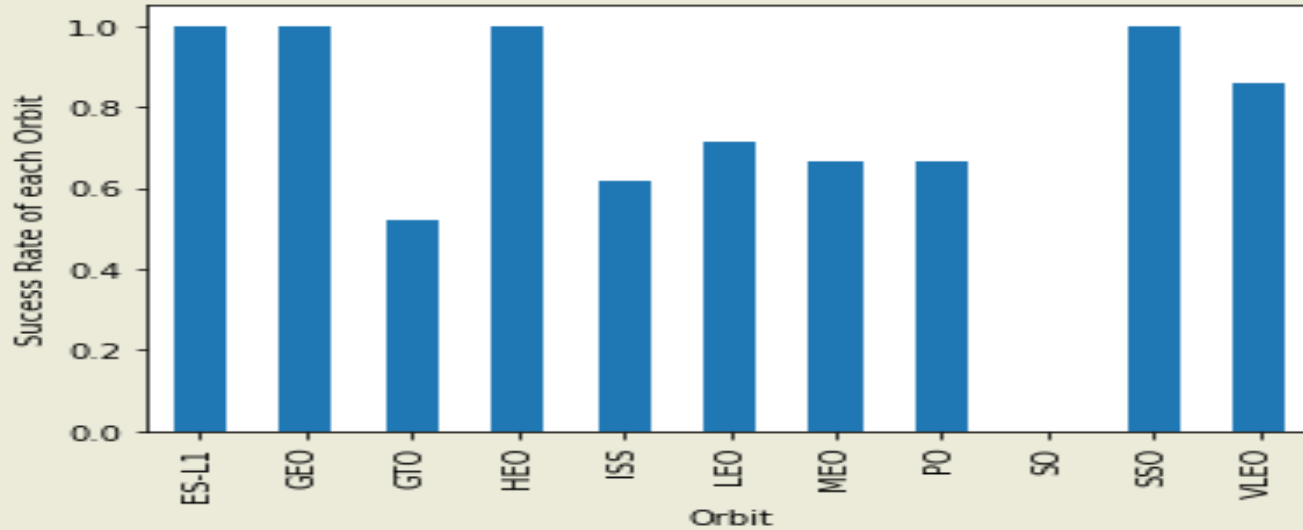
We observe that, for each site, the success rate is increasing.

Payload & Launch Site



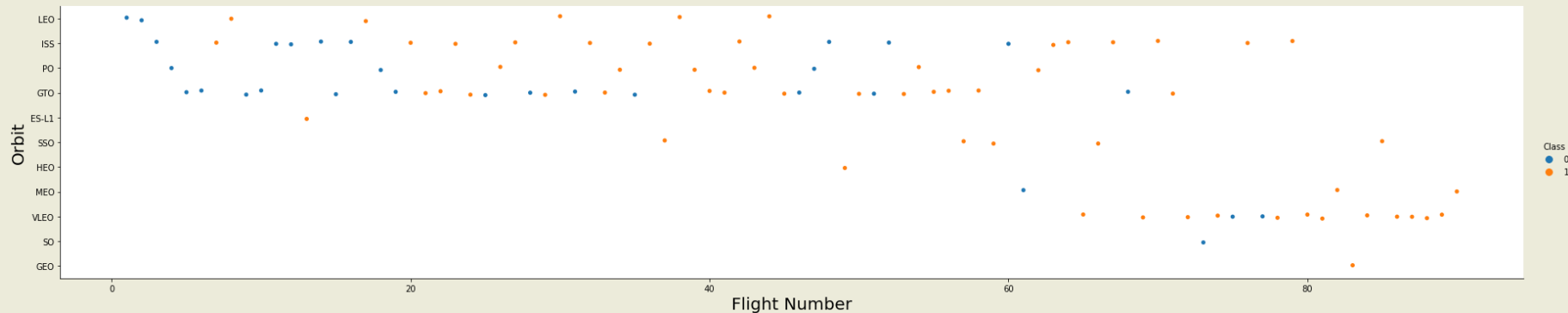
Depending on the launch site, a heavier payload may be a consideration for a successful landing. On the other hand, a too heavy payload can make a landing fail.

Success Rate & Orbit Type



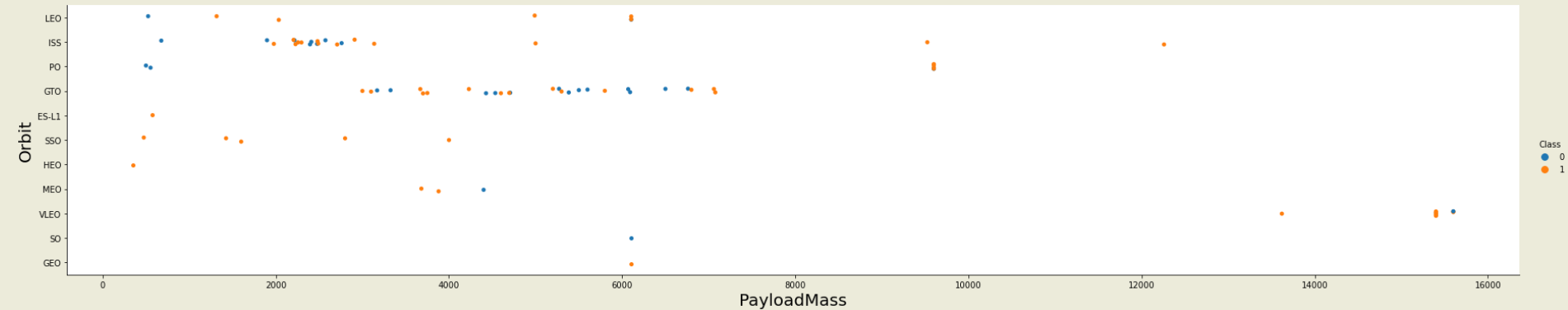
With this plot, we can see success rate for different orbit types. We note that ES-L1, GEO, HEO, SSO have the best success rate.

Flight Number & Orbit Type



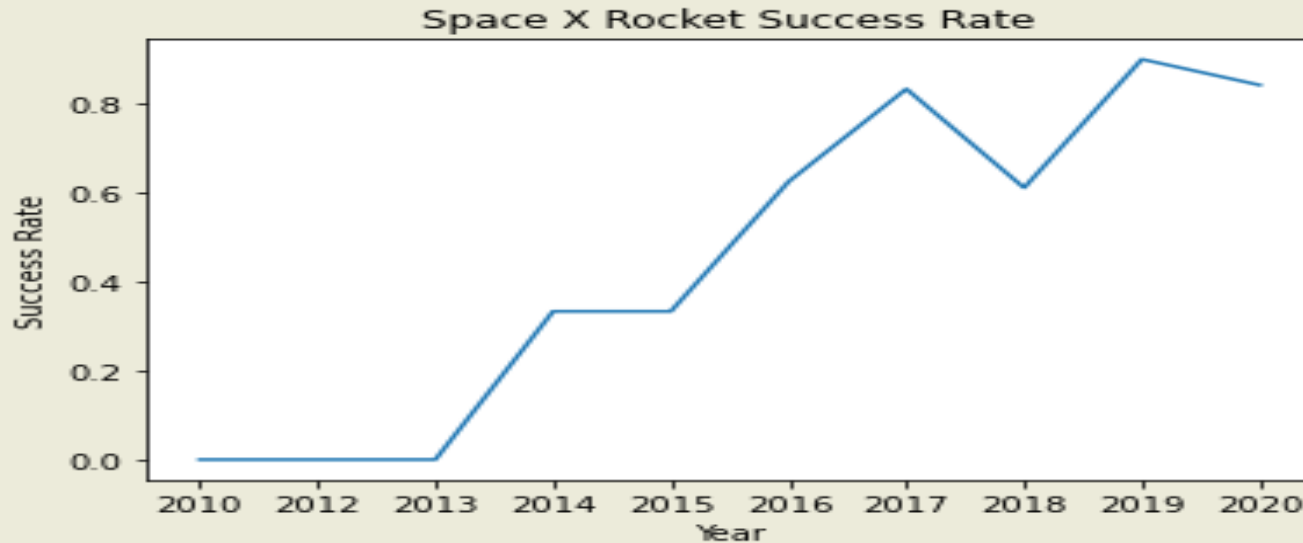
We notice that the success rate increases with the number of flights for the LEO orbit. For some orbits like GTO, there is no relation between the success rate and the number of flights. But we can suppose that the high success rate of some orbits like SSO or HEO is due to the knowledge learned during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

Payload & Orbit Type



Payload weight can significantly affect launch success rates in specific orbits. For instance, heavier payloads tend to enhance success rates in LEO orbits, while reducing payload weight improves launch success for GTO orbits.

Launch Success Yearly Trend



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

All Launch Site Names

SQL Query

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

Results

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

Explanation

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

Launch Site Names Begin with 'CCA'

SQL Query

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

Results

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)

Explanation

The WHERE clause combined with the LIKE operator filters launch sites to include only those containing the substring "CCA." The LIMIT 5 statement then limits the output to 5 records from the filtered results.

Total Payload Mass

SQL Query

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

Results

SUM("PAYLOAD_MASS__KG_")
45596

Explanation

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

Average Payload Mass by F9

SQL Query

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

Results

```
AVG("PAYLOAD_MASS_KG_")  
2534.6666666666665
```

Explanation

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

First Successful Ground Landing Date

SQL Query

```
SELECT MIN(`DATE`) FROM SPACE171 WHERE `landing_outcome` LIKE 'successful'
```

Results

MIN("DATE")

01-05-2017

Explanation

This query retrieves the oldest successful landing. The `WHERE` clause filters the dataset to include only successful landings, and the `MIN` function identifies the record with the earliest date.

Successful Drone Ship Landing with Payload between 4000 and 6000

SQL Query

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

Results

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

Explanation

This query retrieves the booster version for landings that were successful and had a payload mass between 4000 and 6000 kg. The WHERE and AND clauses are used to filter the dataset accordingly.

Total Number of Successful and Failure Mission Outcomes

SQL Query

```
%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
```

Results

SUCCESS	FAILURE
100	1

Explanation

The first SELECT statement displays the results of the subqueries. The first subquery counts the number of successful missions, while the second counts the unsuccessful ones. The WHERE clause with the LIKE operator filters the mission outcomes, and the COUNT function tallies the filtered records.

Boosters Carried Maximum Payload

SQL Query

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

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

Explanation

We used a subquery to filter the data by selecting the heaviest payload mass using the MAX function. The main query then utilizes these results to return distinct booster versions associated with this maximum payload mass.

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'
```

Results

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

Explanation

This query retrieves the month, booster version, and launch site for landings that were unsuccessful in 2015. The `SUBSTR` function extracts parts of the date, with `SUBSTR(DATE, 4, 2)` showing the month and `SUBSTR(DATE, 7, 4)` showing the year.

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

SQL Query

```
%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 ;
```

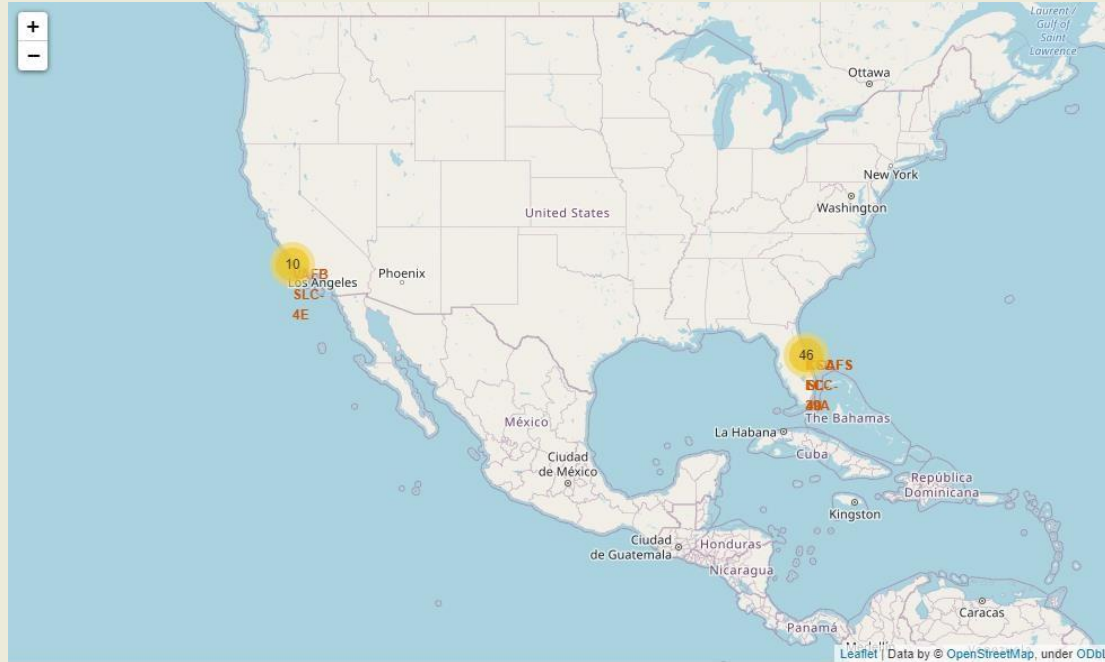
Results

Landing _Outcome	COUNT("LANDING _OUTCOME")
Success	20
Success (drone ship)	8
Success (ground pad)	6

Explanation:

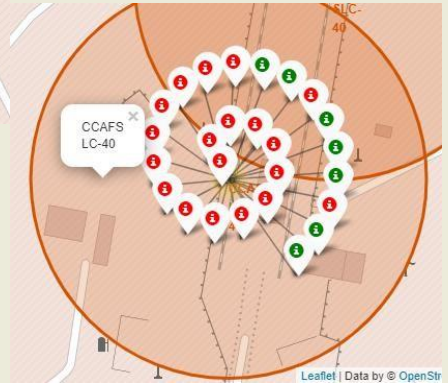
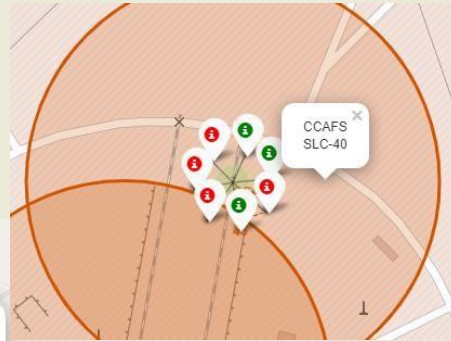
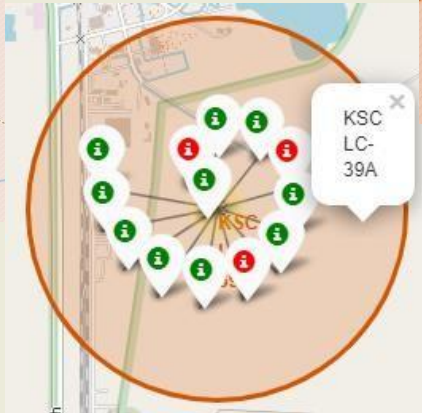
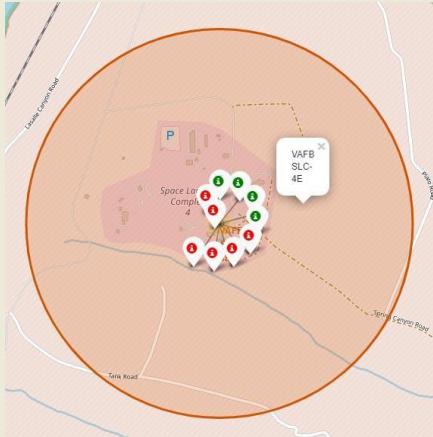
This query retrieves the count of landing outcomes where the mission was successful, with dates ranging from 04/06/2010 to 20/03/2017. Results are grouped by landing outcome using the GROUP BY clause, and sorted in descending order by count with the ORDER BY COUNT DESC statement.

Folium map – Ground stations



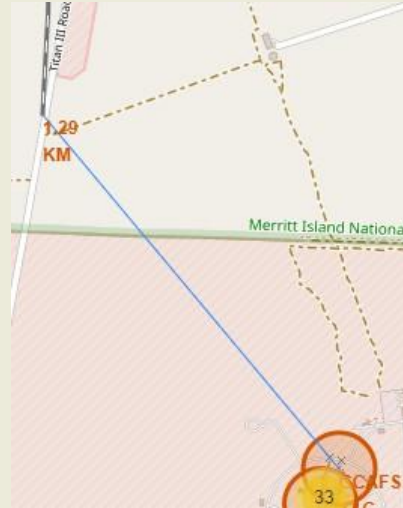
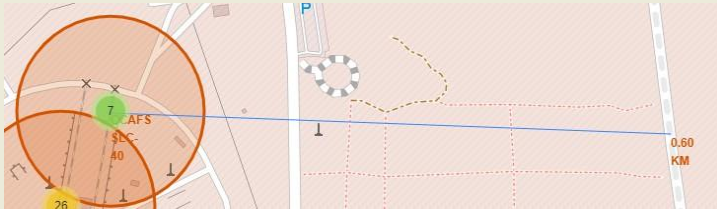
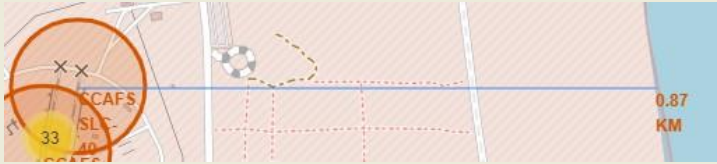
We see that Space X launch sites are located on the coast of the United States

Folium map – Color Labeled Markers



Green marker represents successful launches. Red marker represents unsuccessful launches. We note that KSC LC-39A has a higher launch success rate.

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

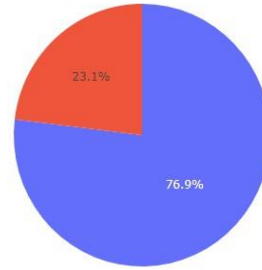
Total Success Launches by Site



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

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

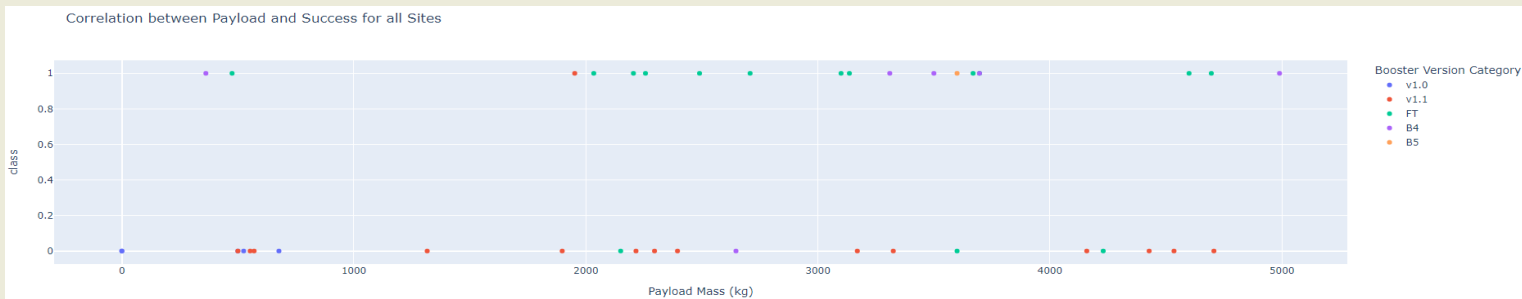
Total Success Launches for Site KSC LC-39A



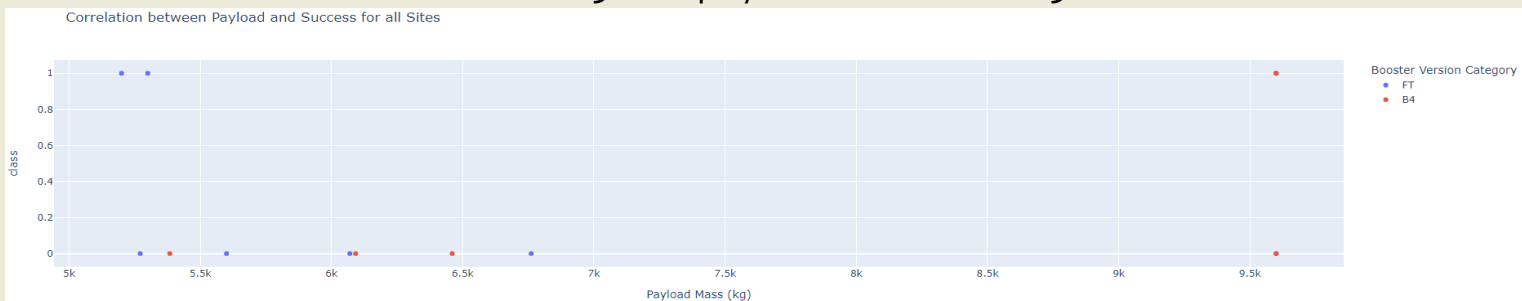
■ 1
■ 0

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

Dashboard – Payload mass & Outcome for all sites with different payload mass selected



Low weighted payload (0 – 5000 kg)

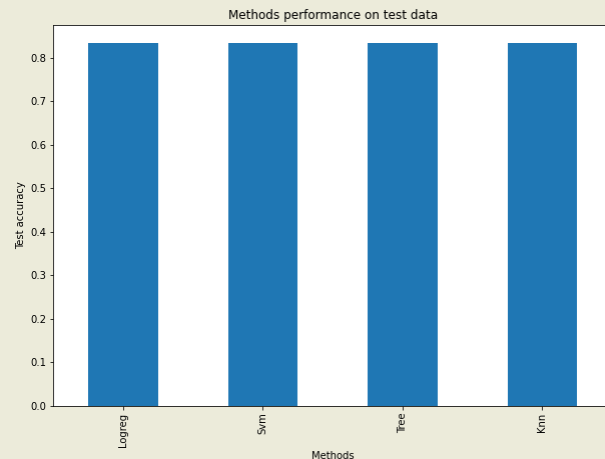
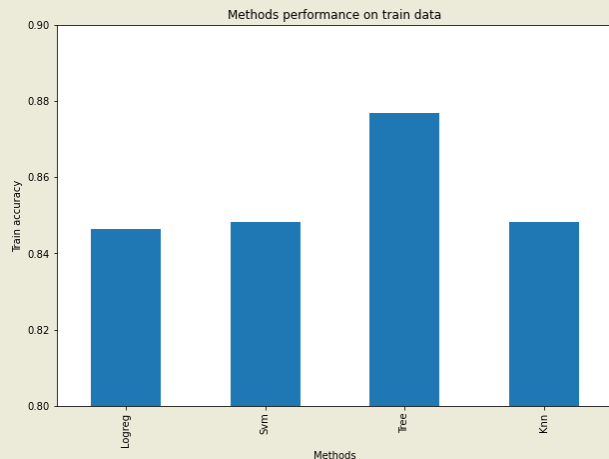


Heavy weighted payload (5000 – 10000 kg)

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



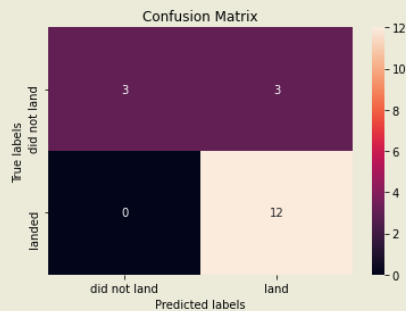
For accuracy test, all methods performed similar. We could get more test data to decide between them. But if we really need to choose one right now, we would take the decision tree.

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

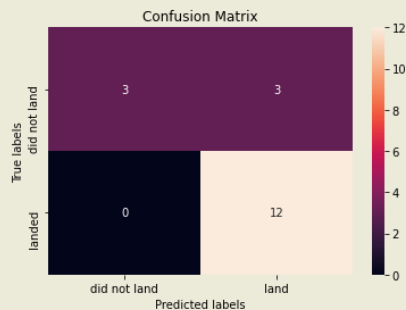


Decision Tree

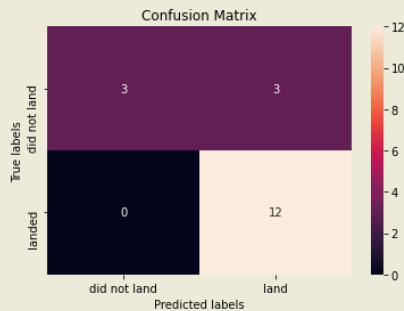


As the test accuracy are all equal, the confusion matrices are also identical. The main problem of these models are false positives.

kNN



SVM



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

- Mission success is affected by factors such as launch site, orbit, and number of prior launches.
- Lighter payloads generally perform better; orbits with high success rates include GEO, HEO, SSO, and ES-L1.
- The exact reasons for the effectiveness of certain launch sites (e.g., KSC LC-39A) are unclear and may require additional data.
- The Decision Tree Algorithm was selected due to its better training accuracy, despite similar test accuracy across models.

O5. Conclu- sion

THANKS!

