

Winning Space Race with Data Science

Daniel Carvalho Pedrini
05/01/2023



Outline

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

Executive Summary

- Summary of methodologies
 - Define and collect the data collection within different ways
 - Define, apply and compare Machine Learning models
 - Create Data Visualization to best view the data
- Summary of all results
 - Define the best model applied
 - Create the Visualizations to have best decision making

Introduction

- Project background and context

For this project, we will work with the SpaceX company to predict the Falcon 9 to check the success of the first step landing. It is important for the matter of costs.

- Problems you want to find answers

1. What is the factors behind the failure of landing?
2. Will rockets land successfully?
3. What is the accuracy of successful landing?

Section 1

Methodology

Methodology

Executive Summary

- Data collection methodology:
 - We used two methodology for data collection: API and Web Scrapping
- Perform data wrangling
 - We used One-hot-encoding
- Perform exploratory data analysis (EDA) using visualization and SQL
- Perform interactive visual analytics using Folium and Plotly Dash
- Perform predictive analysis using classification models
 - We applied different models to compere and find the best model

Data Collection – SpaceX API



<https://github.com/danielpedrini/capstone/blob/main/jupyter-labs-spacex-data-collection-api.ipynb>

Data Collection – SpaceX API

```
static_url = "https://en.wikipedia.org/w/index.php?title=List_of_Falcon_9_and_Falcon_Heavy_launches&oldid=1027686922"

Next, request the HTML page from the above URL and get a response object

TASK 1: Request the Falcon9 Launch Wiki page from its URL

First, let's perform an HTTP GET method to request the Falcon9 Launch HTML page, as an HTTP response.

# use requests.get() method with the provided static_url
# assign the response to a object
response = requests.get(static_url)
```

```
# Use BeautifulSoup() to create a BeautifulSoup object from a response text content
soup = BeautifulSoup(response.text,'html.parser')
```

```
# Use the find_all function in the BeautifulSoup object, with element type 'table'
# Assign the result to a list called 'html_tables'
html_tables = soup.find_all('table')

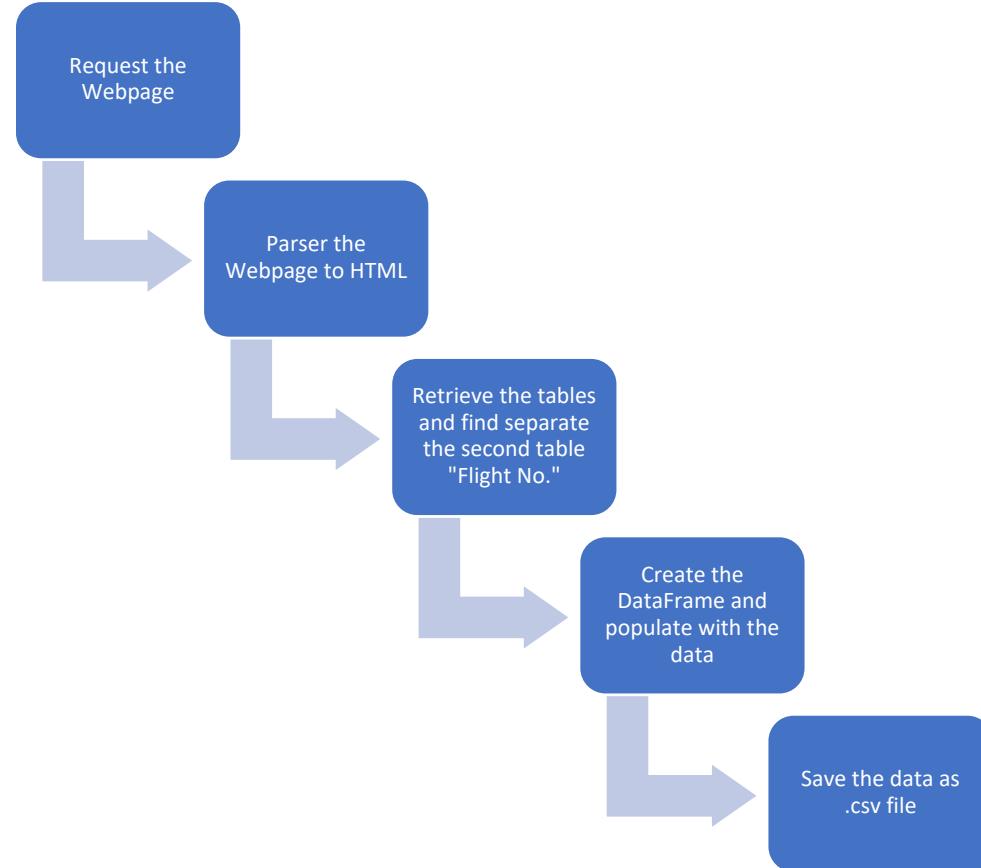
Starting from the third table is our target table contains the actual launch records.

# Let's print the third table and check its content
first_launch_table = html_tables[2]
print(first_launch_table)
```

```
extracted_row = 0
launch_dict = []
for table in html_tables:
    for table_number,table in enumerate(soup.find_all('table'),"visible plainrowheaders collapsible"):
        for rows in table.find_all("tr"):
            if rows.td:
                if rows.th.string:
                    if flight_number==rows.th.string.strip():
                        flagflight_number=True
                else:
                    if flag:
                        extracted_row+=1
                        extracted = rows
                        date = extracted.td[0].text.strip()
                        time = extracted.td[1].text.strip()
                        bvbooster_version=rows[1]
                        launch_site = rows[2].a.string
                        payload = rows[3].a.string
                        payload_mass = rows[4].a.get_attribute('data-mass')
                        orbit = rows[5].a.string
                        if rows[6].a:
                            customer = rows[6].a.string
                        else:
                            customer = None
                        launch_outcome = list(rows[7].strings)[8]
                        booster_landing = landing_status(rows[8])
                        launch_dict.append({
                            'Flight No.': flight_number,
                            'Date': date,
                            'Time': time,
                            'Booster Version': bvbooster_version,
                            'Launch Site': launch_site,
                            'Payload': payload,
                            'Payload mass': payload_mass,
                            'Orbit': orbit,
                            'Customer': customer,
                            'Launch outcome': launch_outcome,
                            'Booster landing': booster_landing
                        })
                        launch_dict.append(launch_outcome)
                    else:
                        flag=False
            rows.find_all('td')[0]
```

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

```
] df.to_csv('sb9cex-Mp-2cl9b6q.csv', index=False)
```



<https://github.com/danielpedrini/capstone/blob/main/jupyter-labs-webscraping.ipynb>

Data Wrangling

Identify missing values and the percentage of missing values

Identify the data type of each columns

Check the unique values of Launch Site, Orbit and Outcomes and the frequency of each values

Create columns Class with failure outcome.

Save the file to format .csv

EDA with Data Visualization

- Charts used in this project.
 - Seaborn catplot: we used this type of graphic to plot how the data behavior between different categories in each class
 - Seaborn barplot: We used this plot to check the mean values of class in each orbit
 - Seaborn lineplot: We used this plot to see the evolution of mean value of the class peer year

<https://github.com/danielpedrini/capstone/blob/main/jupyter-labs-eda-dataviz.ipynb>

EDA with SQL

- **Query Used to perform EDA:**

- `select distinct(Launch_Site) from SPACEXTBL`
- `select * from SPACEXTBL where Launch_Site like 'CCA%' limit 5`
- `select sum(PAYLOAD_MASS__KG_) from SPACEXTBL`
- `select avg(PAYLOAD_MASS__KG_) from SPACEXTBL where Booster_Version like 'F9 v1.1%`
- `select min(date) from SPACEXTBL where "Landing _Outcome" = 'Success'`
- `select distinct(Booster_Version) from SPACEXTBL where "Landing _Outcome" = 'Success (drone ship)' and PAYLOAD_MASS__KG_ between 4000 and 6000`
- `select count(Mission_Outcome) from SPACEXTBL where Mission_Outcome like '%Success'`
- `select count(Mission_Outcome) from SPACEXTBL where Mission_Outcome not like '%Success'`
- `select Booster_Version from SPACEXTBL where PAYLOAD_MASS__KG_ = (select max(PAYLOAD_MASS__KG_) from SPACEXTBL)`
- `elect substr(Date, 4,2) as month, "Landing _Outcome", Booster_Version, launch_site, substr(Date, 7,4) from SPACEXTBL where "Landing _Outcome" = 'Failure (drone ship)' and substr(Date, 7,4) = '2015'`

Build an Interactive Map with Folium

- folium.Marker() used to create marks on map
- folium.Circle() used to create a circles above ther markers
- folium.Icon() used to create an icon on the map
- folion.PolyLine() used to create polynomial line between points
- folium.plugins.AntPath() used to create animated line between points
- markerCluster used to simplifiy the map

https://github.com/danielpedrini/capstone/blob/main/lab_jupyter_launch_site_location.ipynb

Build a Dashboard with Plotly Dash

- Dash and html components we were used and almost everything depends on them like graphs, tables, dropbox
- Plotly used to plot graphs
- Pandas used to manipulate the data
- Piechart and scatterchart to present the chart data
- Slider to select the mass range
- Dropdown to select the sites

<https://github.com/danielpedrini/capstone/blob/main/dashboard.py>

Predictive Analysis (Classification)

Separate the Features and the Target in two DataFrame

Transform the Feature to best predict the model using StandardScaler

Split the data into X_train, X_test, Y_train and Y_test

Select four models to analyze: Logistic Regression, SVM, Decision Tree and KNN

To find the best tuning, we used GridSearch to set the parameters and test all the parameters defined

Tested all the models defined with some parameters

Check the results with: Accuracy, R2_score for train and test Values

Check the Confusion Matrix for each best parameters result on the test

Compare best parameters result between the models

Results

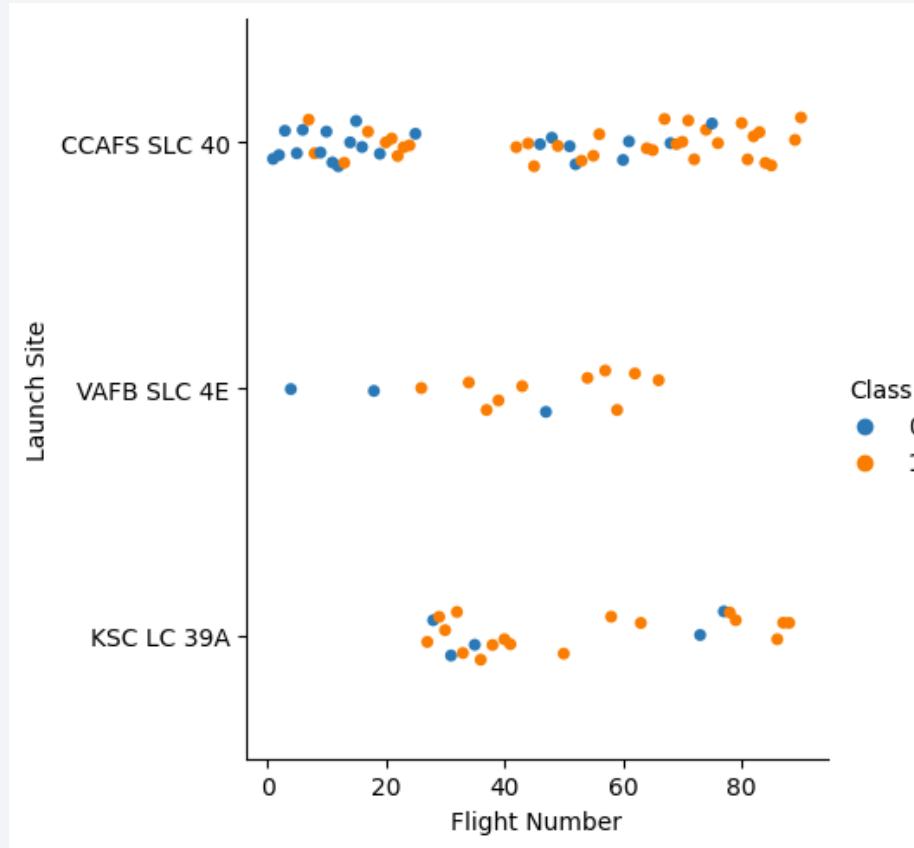
- The site with highest score was KSC LC-39A
- The Payload of 0 to 5000kg was more diverse than 6000kg to 10000kg
- Decision Tree was the best optimal model with the accuracy of 0.8732

The background of the slide features a complex, abstract digital visualization. It consists of numerous thin, glowing lines that create a sense of depth and motion. The lines are primarily blue and red, with some green and purple highlights. They form a grid-like structure that curves and twists across the frame, resembling a three-dimensional space or a network of data points. The overall effect is futuristic and dynamic.

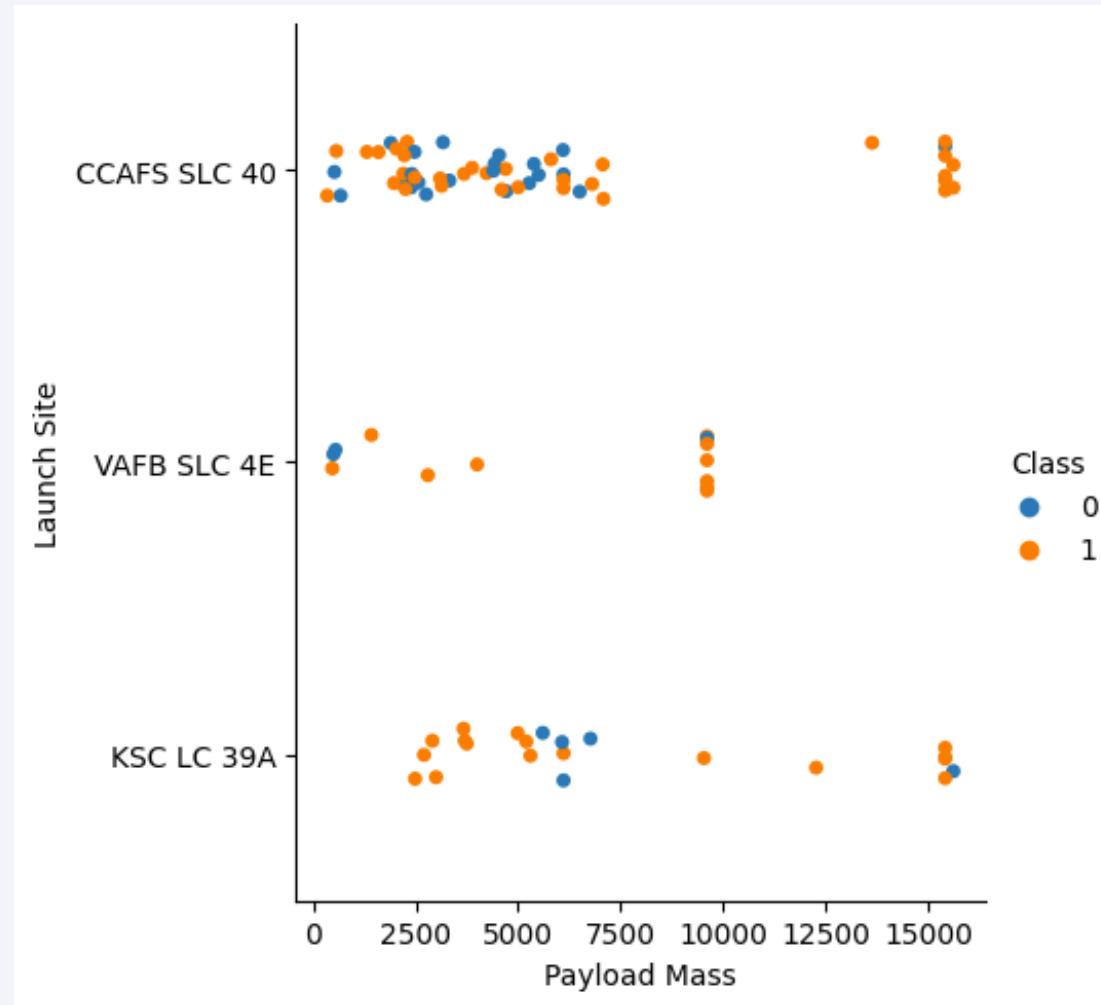
Section 2

Insights drawn from EDA

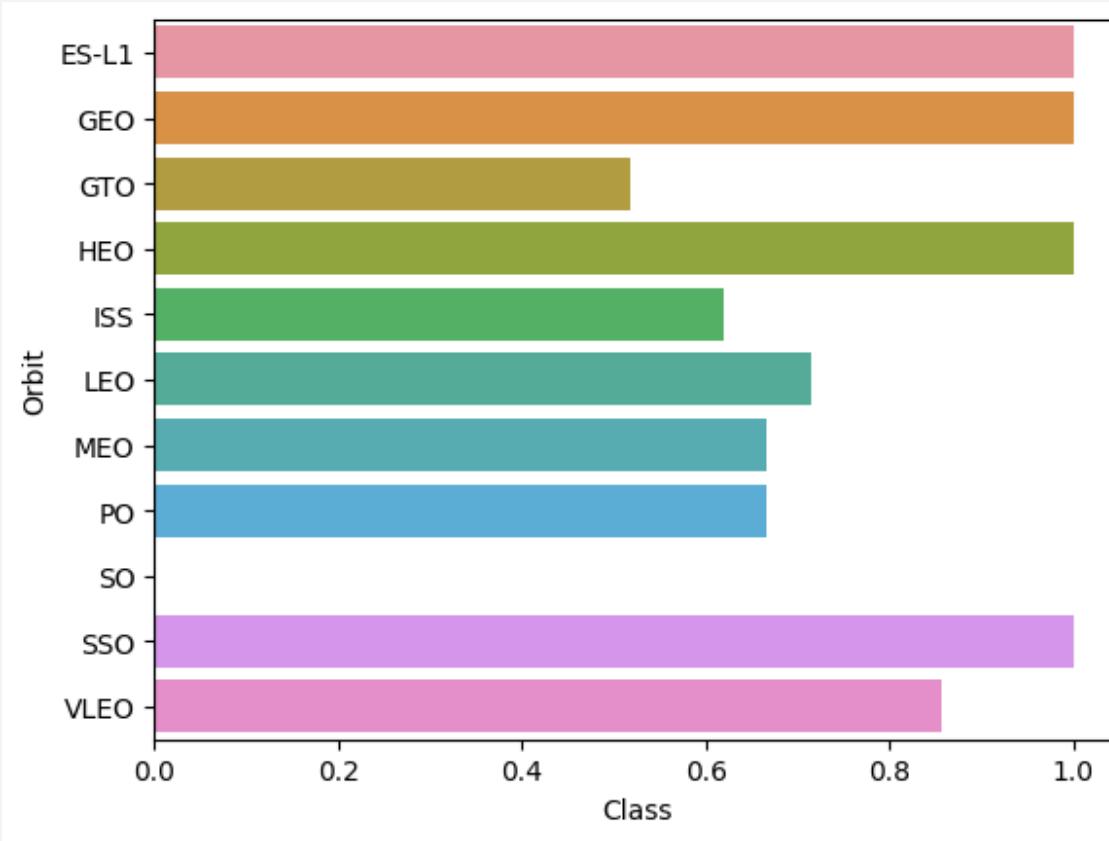
Flight Number vs. Launch Site



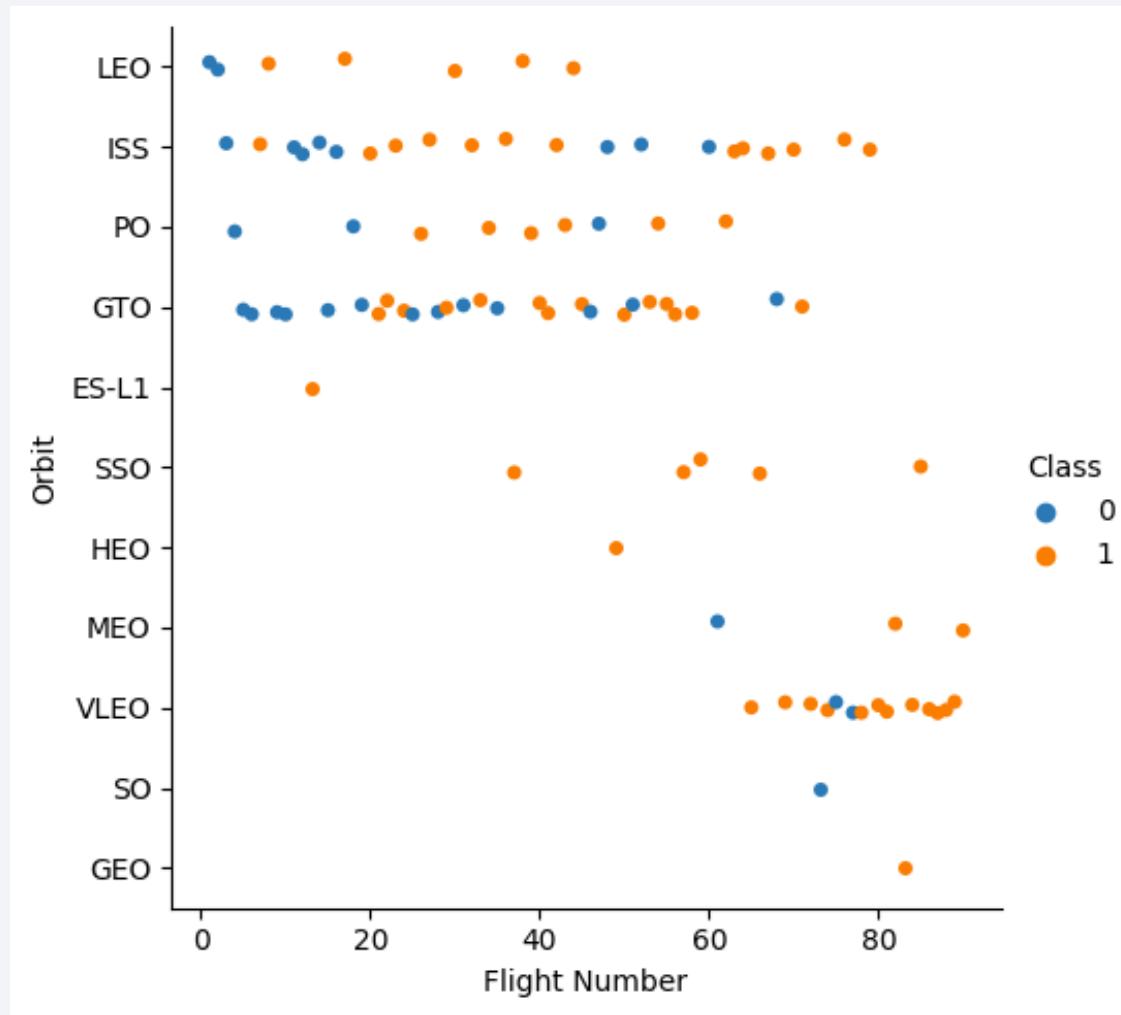
Payload vs. Launch Site



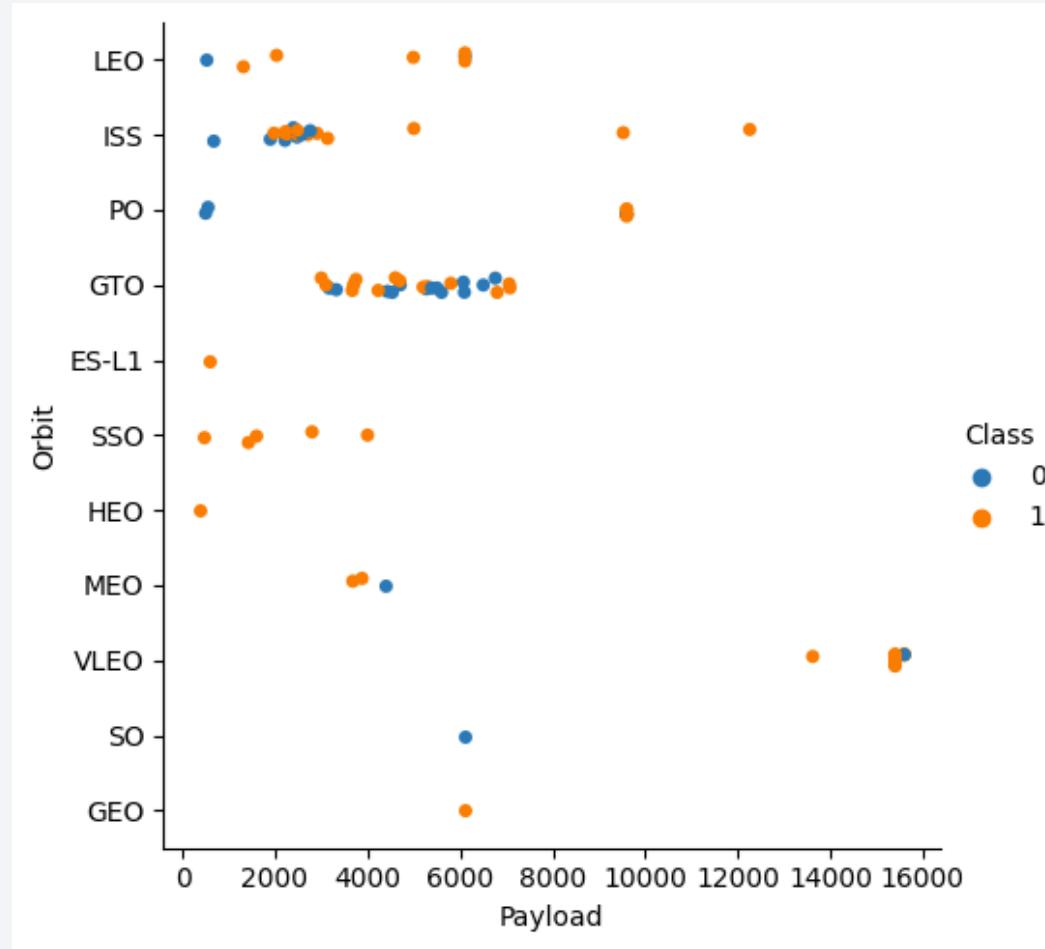
Success Rate vs. Orbit Type



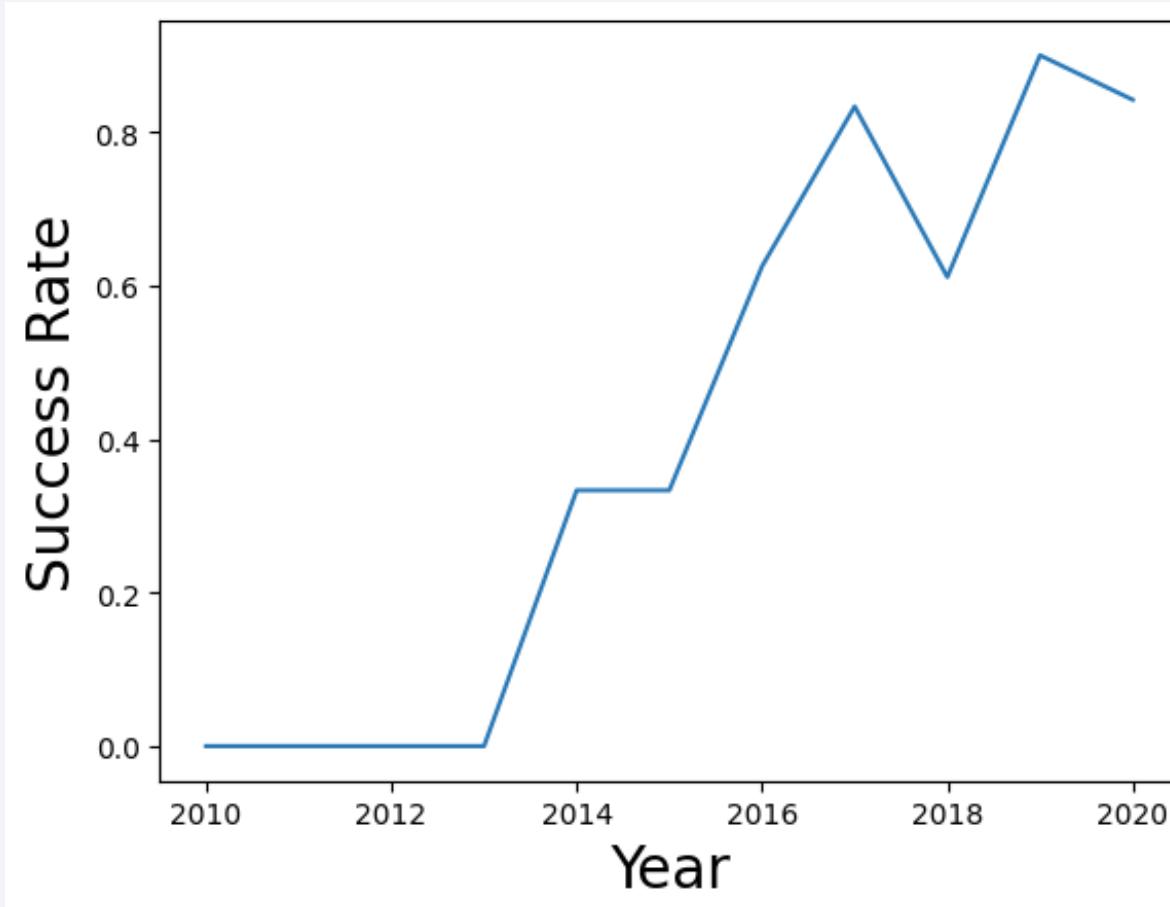
Flight Number vs. Orbit Type



Payload vs. Orbit Type



Launch Success Yearly Trend



All Launch Site Names

```
%sql SELECT distinct(Launch_Site) FROM SPACEXTBL

✓ 0.0s
* sqlite:///my_data1.db
Done.

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

Launch Site Names Begin with 'CCA'

```
%%sql sqlite:///my_data1.db
select launch_site from SPACEXTBL where Launch_Site like 'CCA%' limit 5
✓ 0.0s
```

Done.

Launch_Site

CCAFS LC-40

CCAFS LC-40

CCAFS LC-40

CCAFS LC-40

CCAFS LC-40

Total Payload Mass

```
%sql select sum(PAYLOAD_MASS__KG_) from SPACEXTBL
✓ 0.0s
* sqlite:///my_data1.db
Done.

sum(PAYLOAD_MASS__KG_)
619967
```

Average Payload Mass by F9 v1.1

```
%sql select avg(PAYLOAD_MASS__KG_) from SPACEXTBL where Booster_Version like 'F9 v1.1'
✓ 0.0s
* sqlite:///my\_data1.db
Done.

avg(PAYLOAD_MASS__KG_)
2534.6666666666665
```

First Successful Ground Landing Date

```
%sql select min(date) from SPACEXTBL where "Landing _Outcome" = 'Success'  
✓ 0.0s  
  
* sqlite:///my\_data1.db  
Done.  
  
min(date)  
02-03-2019
```

Successful Drone Ship Landing with Payload between 4000 and 6000

```
%sql select distinct(Booster_Version) from SPACEXTBL where "Landing _Outcome" = 'Success (drone ship)' and PAYLOAD_MASS__KG_ between 4000 and 6000
✓ 0.0s
* sqlite:///my_data1.db
Done.

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

Total Number of Successful and Failure Mission Outcomes

```
%sql select count(Mission_Outcome) from SPACEXTBL where Mission_Outcome like '%Success';  
✓ 0.0s  
* sqlite:///my_data1.db  
Done.  
  
count(Mission_Outcome)  
98  
  
%sql select count(Mission_Outcome) from SPACEXTBL where Mission_Outcome not like '%Success';  
✓ 0.0s  
* sqlite:///my_data1.db  
Done.  
  
count(Mission_Outcome)  
3
```

Boosters Carried Maximum Payload

```
%sql select Booster_Version from SPACEXTBL where PAYLOAD_MASS__KG_ = (select max(PAYLOAD_MASS__KG_) from SPACEXTBL)
✓ 0.0s
* sqlite:///my_data1.db
Done.

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

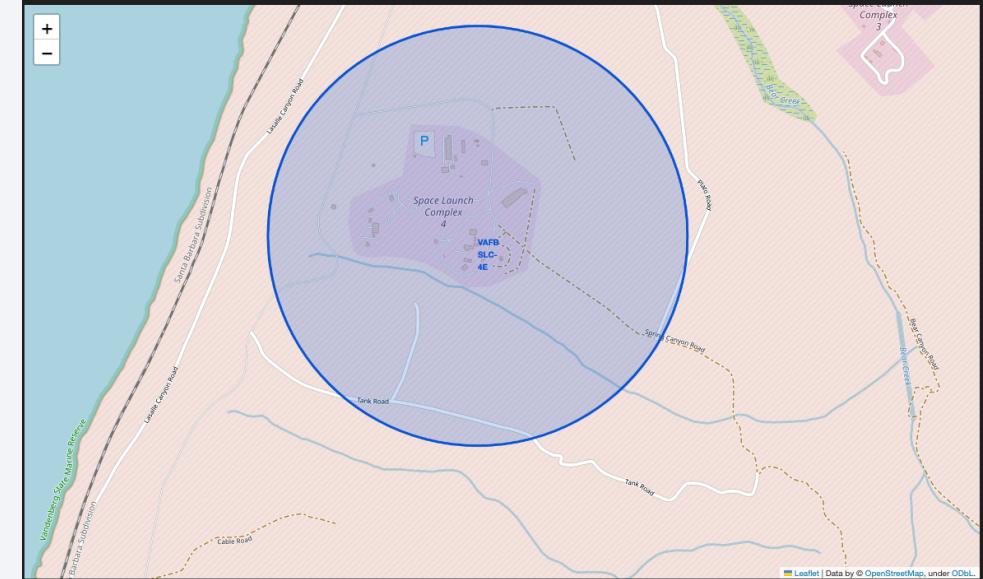
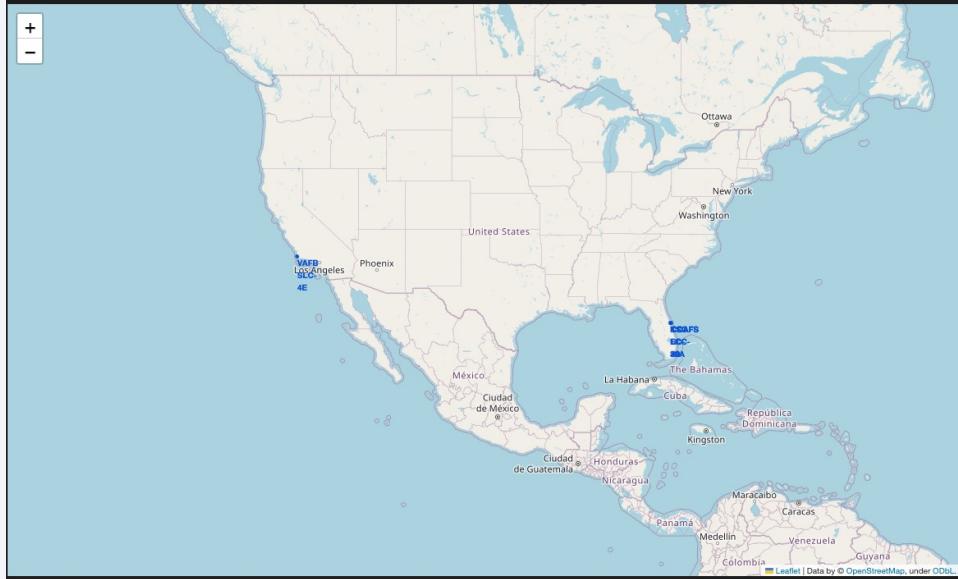
```
_Outcome", Booster_Version, launch_site, substr(Date, 7,4) from SPACEXTBL where "Landing _Outcome" = 'Failure (drone ship)' and substr(Date, 7,4) = '2015'  
4]   ✓  0.0s                                         Python  
* sqlite:///my_data1.db  
Done.  
  
month  Landing _Outcome  Booster_Version  Launch_Site  substr(Date, 7,4)  
01    Failure (drone ship)  F9 v1.1 B1012  CCAFS LC-40      2015  
04    Failure (drone ship)  F9 v1.1 B1015  CCAFS LC-40      2015
```

The background of the slide is a photograph taken from space at night. It shows the curvature of the Earth against a dark blue-black void of space. City lights are visible as numerous small white and yellow dots, primarily concentrated in the lower right quadrant where the United States appears. In the upper right, the green and yellow glow of the aurora borealis is visible. The overall atmosphere is mysterious and scientific.

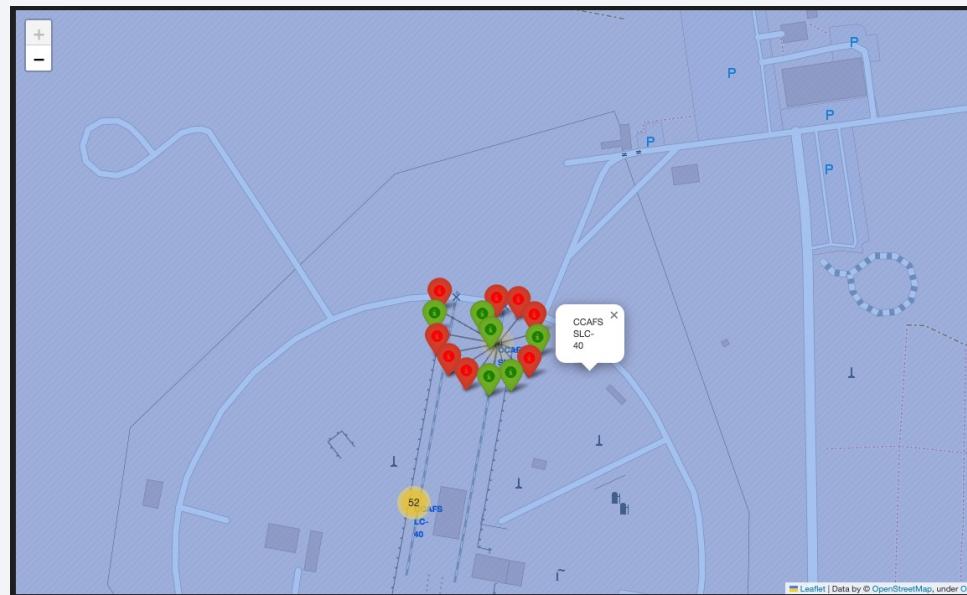
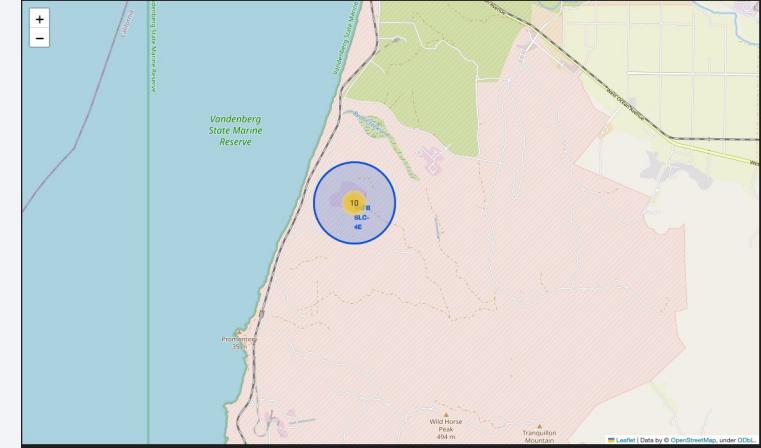
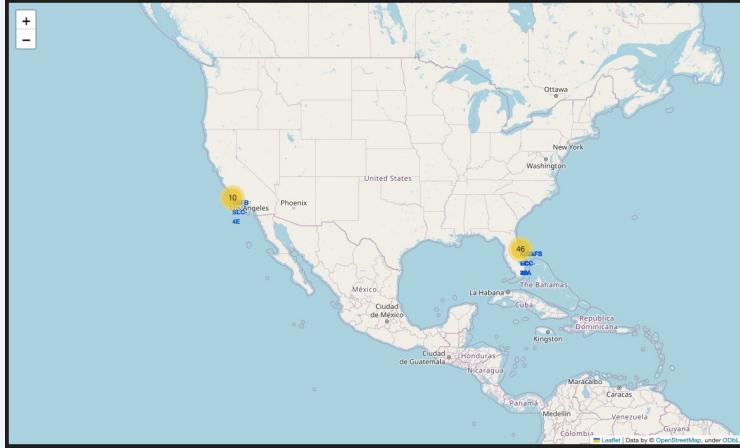
Section 3

Launch Sites Proximities Analysis

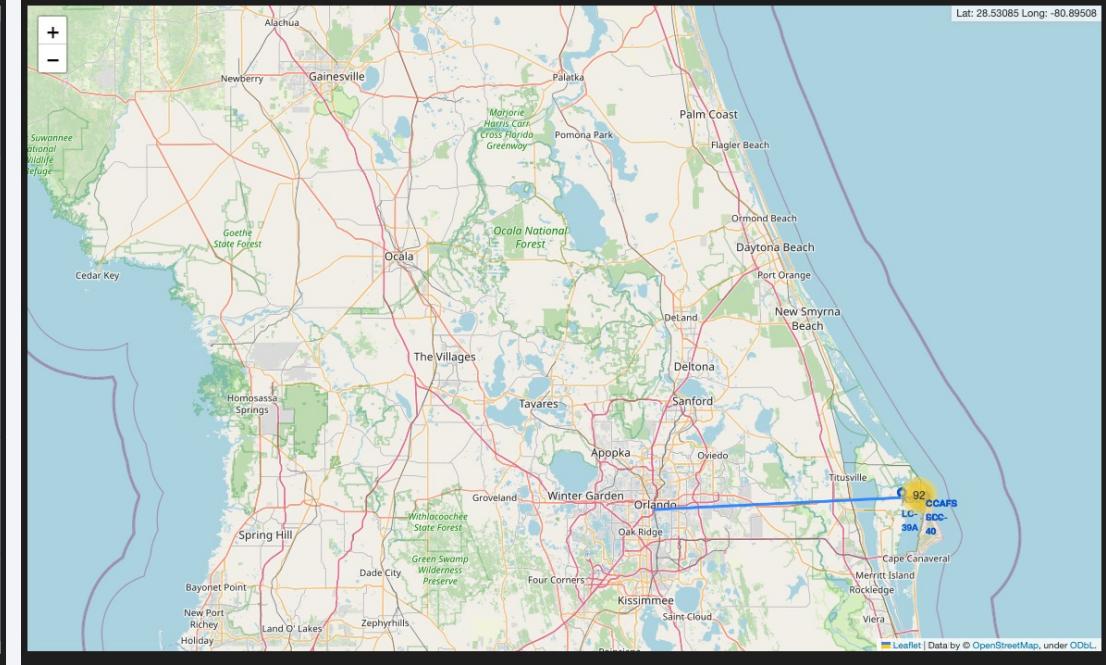
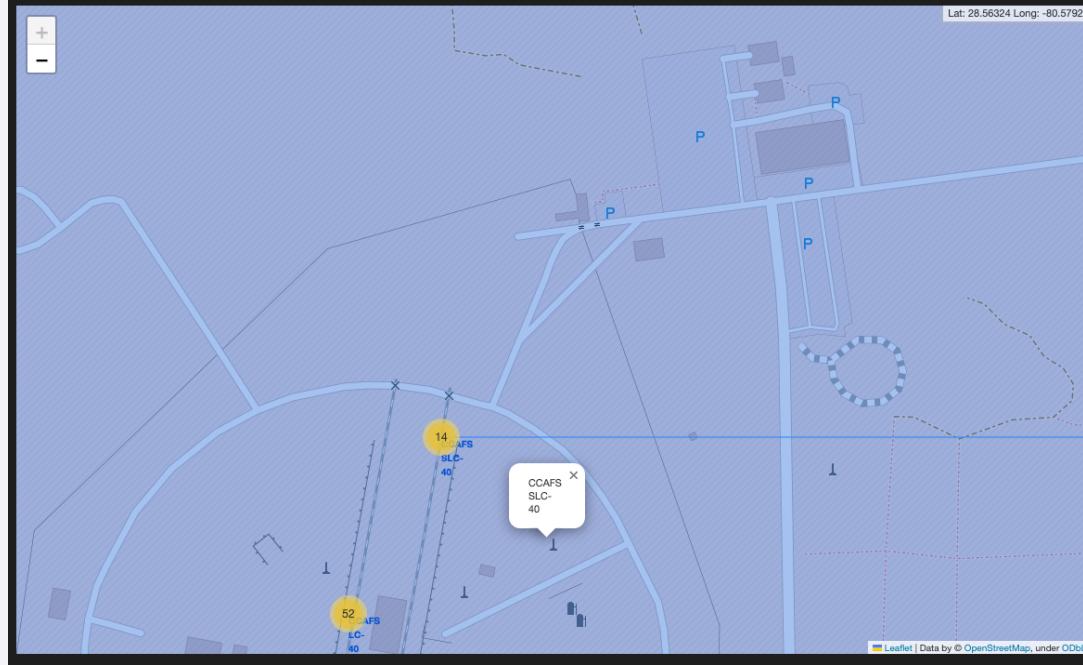
Map One – Marks and Circle



Map 2 – Cluster and Success /Failure

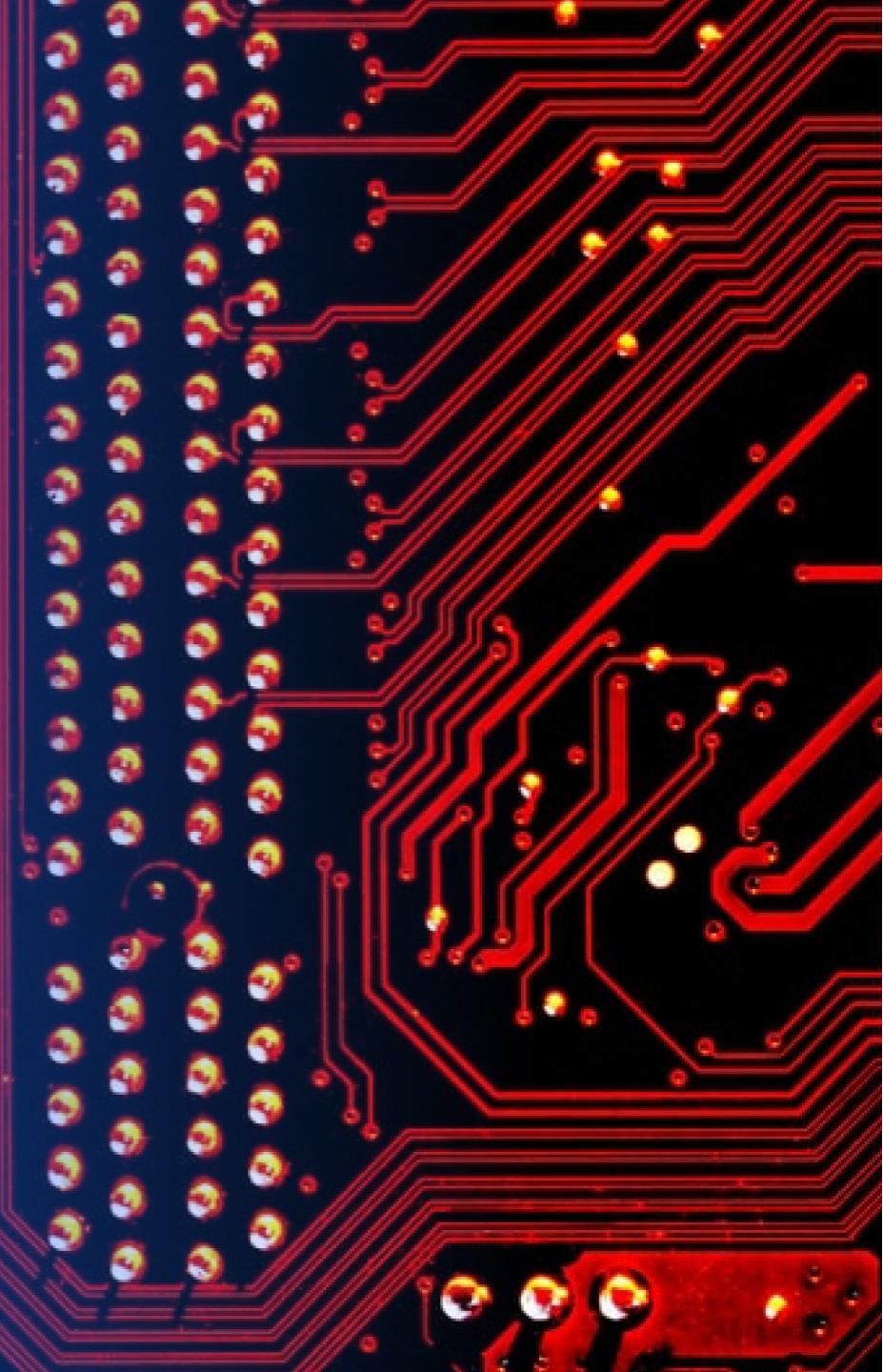


Map 3 - Line

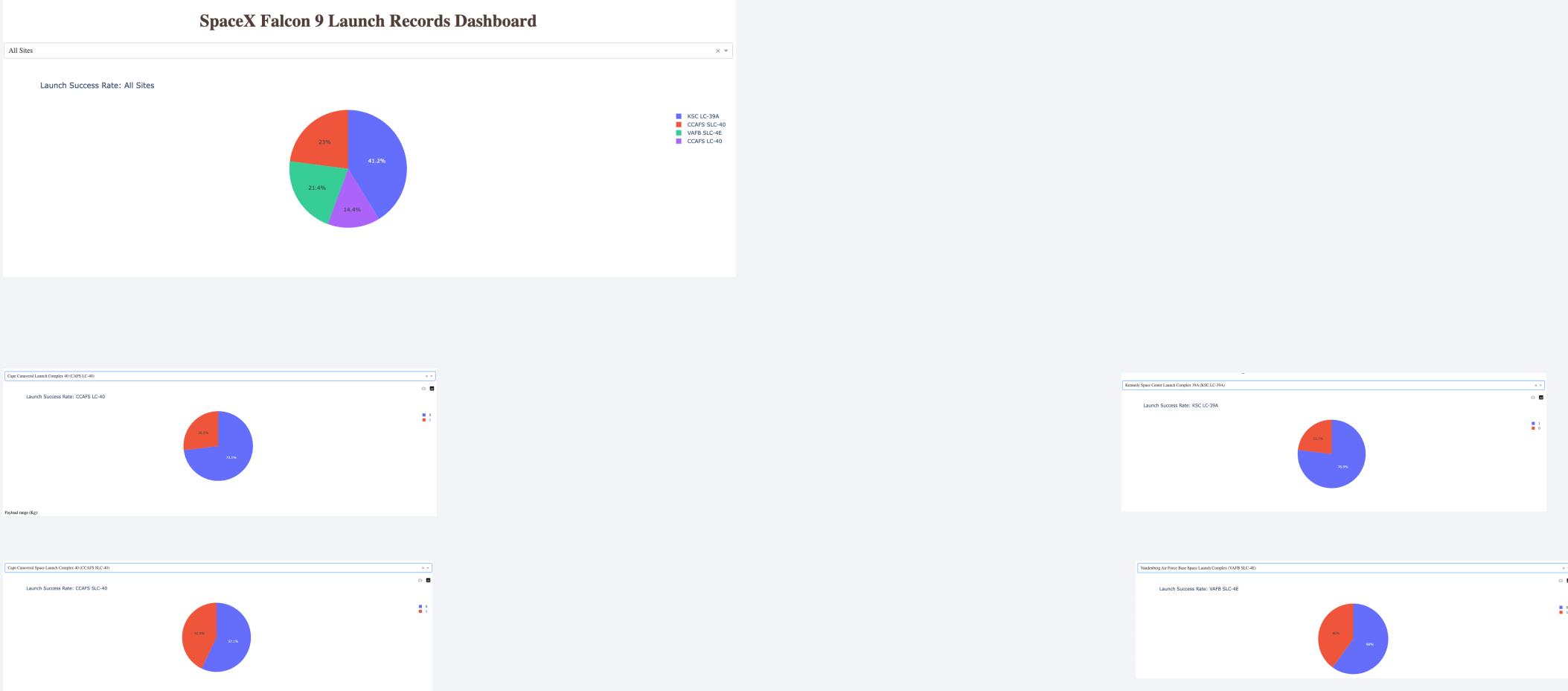


Section 4

Build a Dashboard with Plotly Dash



Dashboard



Dashboard highest success

SpaceX Falcon 9 Launch Records Dashboard

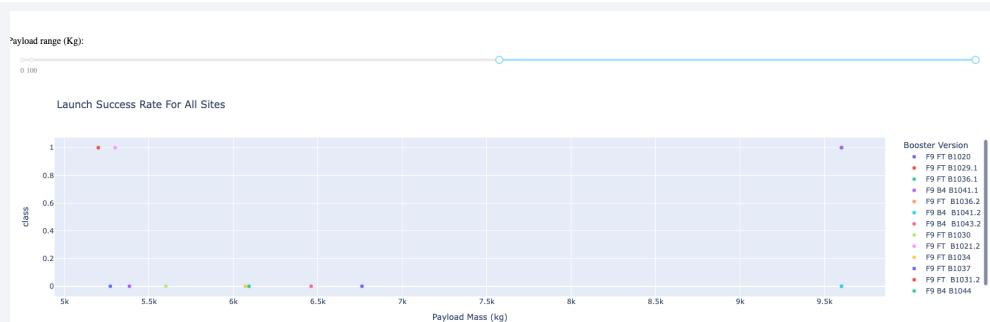
Kennedy Space Center Launch Complex 39A (KSC LC-39A) ✖ ▾

Launch Success Rate: KSC LC-39A



Payload range (Kg):

<Dashboard Screenshot 3>

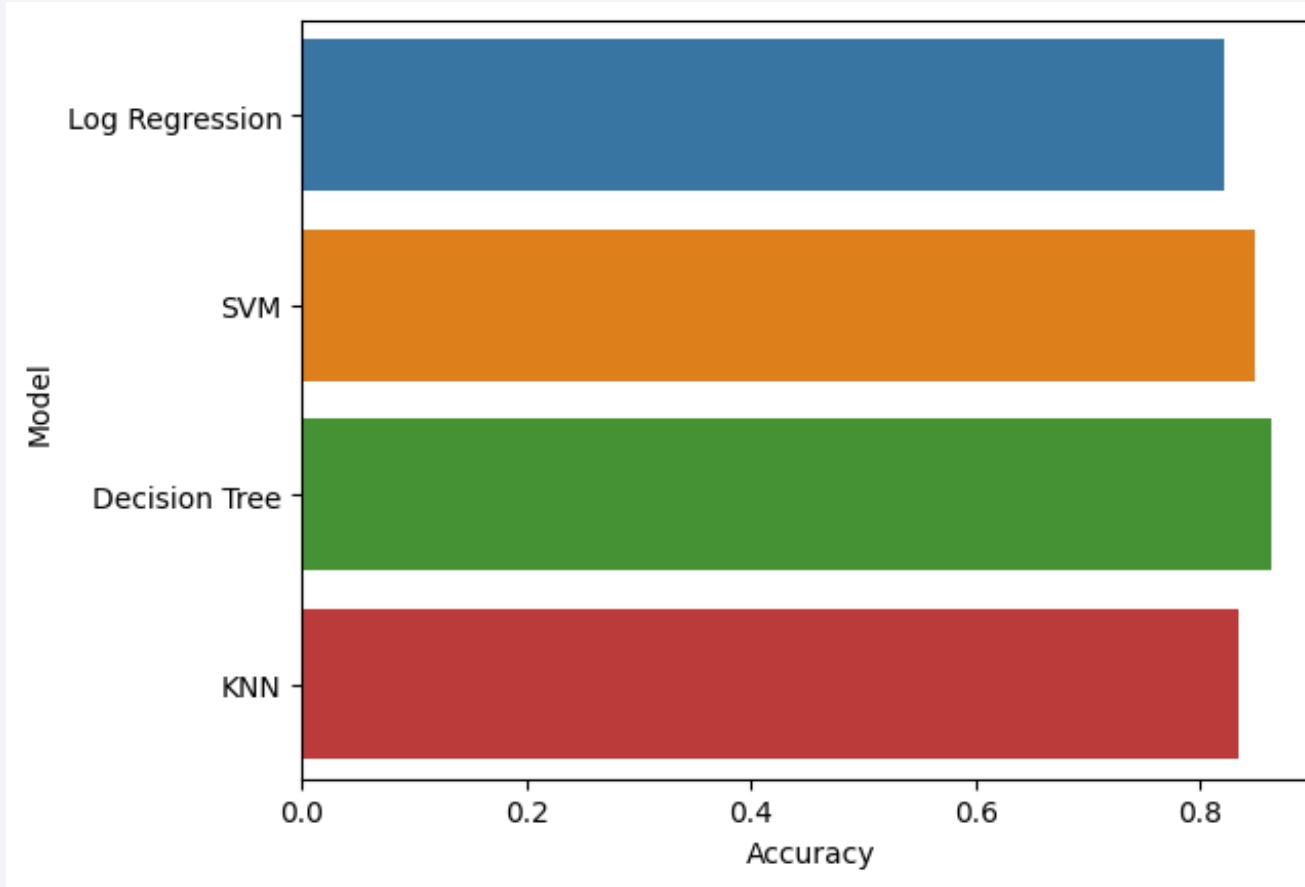


The background of the slide features a dynamic, abstract design. It consists of several thick, curved lines that transition from a bright yellow at the top right to a deep blue at the bottom left. These lines create a sense of motion and depth, resembling a tunnel or a stylized road. The overall effect is modern and professional.

Section 5

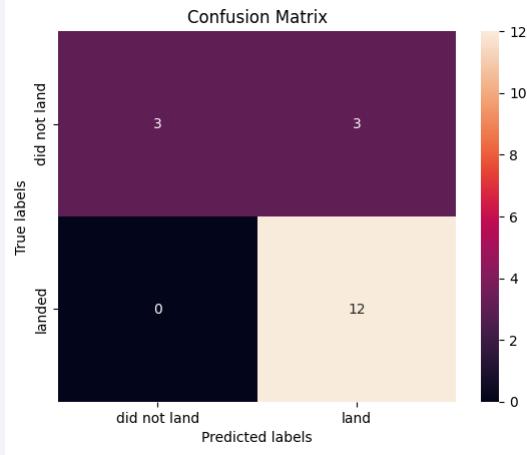
Predictive Analysis (Classification)

Classification Accuracy

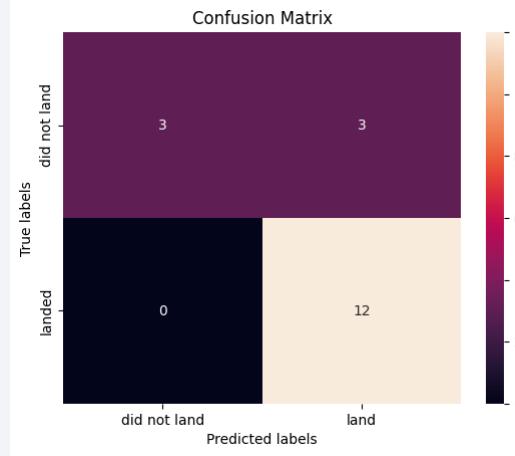


Confusion Matrix

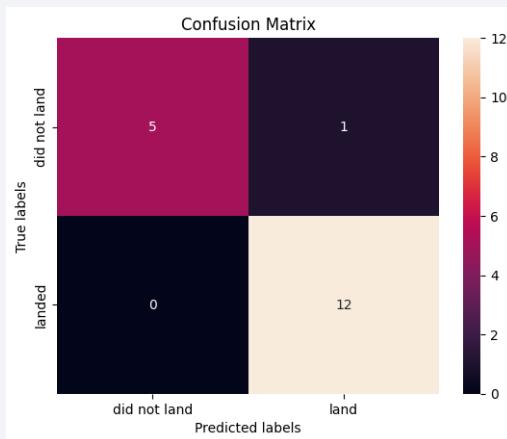
Logistic Regression



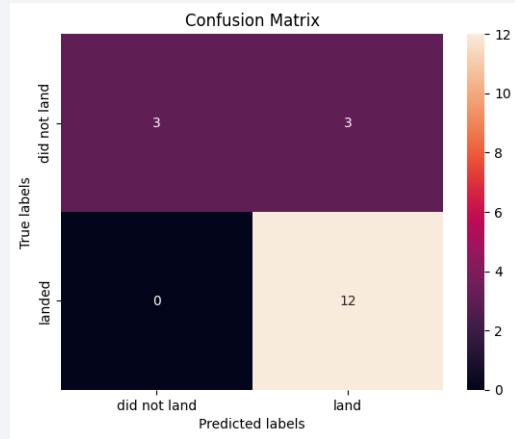
SVM



Decision Tree



KNN



Conclusions

- Decision Tree has the best Accuracy
- Comparing the R2_score in Test data, Decision tree has the best result to, with 0.94
- Using the GridSearch, we could find the best hyperparameter for each model, and with the results compare each one to decide the best model

Appendix

- Access to the github project

<https://github.com/danielpedrini/capstone>

Thank you!

