

# Winning Space Race with Data Science

Nikki de Vries November 6, 2024



### **Table of Contents**



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

### **Executive Summary**

- This project aimed to investigate and find drivers that influenced a successful first-stage landing of the SpaceX Falcon 9 rocket.
- To get to this goal, SpaceX mission data was obtained through a public API and launch data from web scrapping Wikipedia. The data was then cleaned, processed, and transformed for Exploratory Data Analysis (EDA) using Python and SQL. Visualizations such as graphs, maps, and interactive dashboards to find drivers.
- Machine learning classification models, logistic regression, support vector machines, decision trees, and K-nearest-neighbors were analyzed for their training and test accuracy for predicting launch success.
- Based on the results from the analysis:
  - The factors that influenced the successful first-stage landing are time, payload mass, and desired orbit
  - Using these parameters, each classification model had an accuracy of 83.3%, predicting the first-stage rocket booster landing successfully. The decision tree model had the best train accuracy.

### Introduction

- 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.
- As a data scientist, if I determine a successful first-stage landing for the Falcon 9 rocket, I can determine the cost of the launch. This information will help SpaceX prevent alternative companies from bidding against a SpaceX launch, allowing SpaceX to continue its Falcon 9 launches and continue innovating.
- To find the answers to predict a successful first-stage launch, I need to know:
  - Which Factors determine if the Falcon 9 Launch will land successfully?
  - Are there interactions between features that influence this success rate?
  - Are there additional insights I can gather from the data that may also benefit cost?



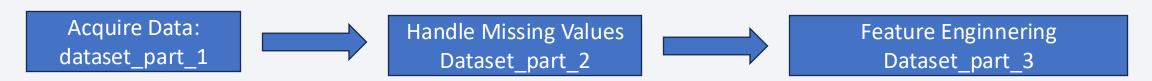
# Methodology

#### **Executive Summary**

- Data collection methodology:
  - Data was collected through the SpaceX API and web scrapping Wikipedia
- Perform data wrangling
  - Missing values were replaced, and a landing outcome feature was created.
- Perform exploratory data analysis (EDA) using visualization and SQL
- Perform interactive visual analytics using Folium and Plotly Dash
- Perform predictive analysis using classification models
  - Normalized data and then performed a train test split for four classification models: logistic regression, support vector machines, decision trees, and k-nearest-neighbors.

#### **Data Collection Overview**

- API
  - Acquired API data through a get request to the SpaceX API.
  - Cleaned the requested data

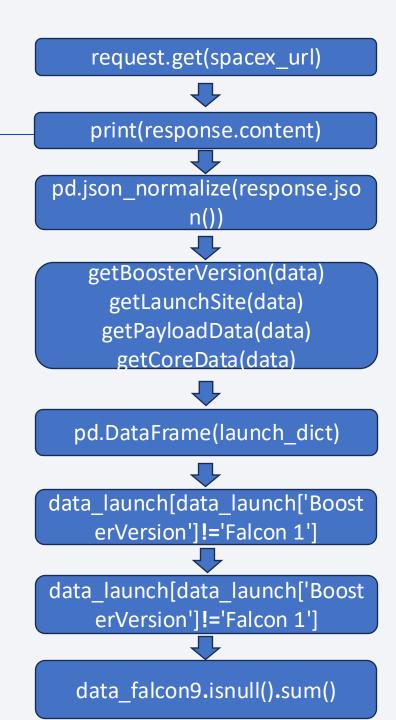


- Web scraping
  - Extracted Falcon 9 launch records HTML table from Wikipedia
  - Parsed the table and converted it into a pandas data frame

# Data Collection – SpaceX API

- 1. Request rocket launch data.
  - Print data to check the content of the response
- 2. Decode the response content as a JSON
  - Normalizing the data to convert the json result into a data frame
- 3. Use API to request and replace column information
  - Take a subset of the data frame to keep the desired features
  - Create request functions for rocket data, launchpad data, payload data, core data
- 4. Create data frame
  - · Combined data into a dictionary to be converted to a dataframe
- 5. Filter the data to only include Falcon 9 launches
- 6. Deal with missing Values
- 7. Save to CSV

Github URL: <a href="https://github.com/NikkideVries/IBM-SpaceX-Capstone/blob/main/1Data\_Collection\_API.ipynb">https://github.com/NikkideVries/IBM-SpaceX-Capstone/blob/main/1Data\_Collection\_API.ipynb</a>

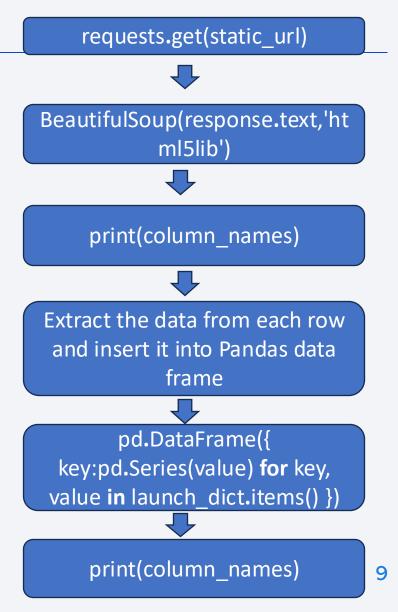


### Data Collection – Web Scraping

- 1. Request the Falcon 9 launch data from the Wiki page URL
- 2. Use BeautifulSoup to parse the content
- 3. Extract the columns and variables from the HTML
- 4. Create data frame by parsing the launch HTML tables
- 5. Export to CSV

Git Hub: <a href="https://github.com/NikkideVries/IBM-SpaceX-">https://github.com/NikkideVries/IBM-SpaceX-</a>

<u>Capstone/blob/main/2Data\_Collection\_with\_Web\_S</u> <u>crapping.ipynb</u>



# **Data Wrangling**

Identify numerical and categorical columns

- Completed some initial exploratory data analysis to see what was happening in the data.
- Figured out from the data what mission outcome was either a failure to land or a successful landing.
- Based on that information, created a new classification variable, "class". O indicates a bad outcome while 1 indicates a good outcome.
  - · Bad outcomes are failures to land
  - Good outcomes are successful landings

Github: <a href="https://github.com/NikkideVries/IBM-SpaceX-Capstone/blob/main/3Data\_Wrangling.ipynb">https://github.com/NikkideVries/IBM-SpaceX-Capstone/blob/main/3Data\_Wrangling.ipynb</a>

Calculate the number of launches on each site



Calculate the number of occurrences of each orbit



Calculate the number of mission outcomes



Encode Outcomes:
Bad outcomes are 0, and good
outcomes are 1



Add new variables to data frame

### **EDA** with Data Visualization

- Exploratory Data Analysis was completed by using Matplotlib visualizations. These visualizations were catplots, scatterplots, bar graphs, and line graphs. Specifically:
- Visualizations were color-coded by class to see successes and failures to land.
- The Features were engineered so categorical variables became numerical to allow for further visualization.

GITHUB: <a href="https://github.com/NikkideVries/IBM-SpaceX-Capstone/blob/main/5EDA">https://github.com/NikkideVries/IBM-SpaceX-Capstone/blob/main/5EDA</a> with Visualizations.ipynb

- Cat plots for:
  - Flight Number and Launch Site
  - Payload Mass and Launch Site
- Bar plot for:
  - Success rate and Orbit Type
- Scatter Plot for:
  - Flight Number and Orbit Type
  - Payload and Orbit Type
- Line Graph for:
  - Year and Success Rate

### **EDA** with SQL

- Displaced the names of the unique launch sites in the space mission. SELECT DISTINCT Launch\_site FROM SPACEXTBL;
- Displayed 5 records where launch sites begins with the string 'CCA'. SELECT \* FROM SPACEXTBL WHERE Launch\_Site LIKE 'CCA%' LIMIT 5;
- Displayed the total payload mass carried by boosters launched by NASA. SELECT SUM(PAYLOAD\_MASS\_\_KG\_) AS TotalPayloadMass FROM SPACEXTBL WHERE Customer = 'NASA (CRS)'
- Display average payload mass carried by booster version F9 v1.1. **SELECT AVG(PAYLOAD\_MASS\_\_KG\_) FROM SPACEXTBL WHERE Booster\_Version = 'F9 v1.1'**
- List the date when the first successful landing outcome in ground pad was acheived. **SELECT MIN(Date) AS EarliestDate FROM SPACEXTBL WHERE Landing\_Outcome LIKE** '%(ground pad)' AND Mission\_Outcome = 'Success';
- List the names of the boosters which have success in drone ship and have payload mass greater than 4000 but less than 6000. SELECT
   Booster\_Version FROM SPACEXTBL WHERE Landing\_Outcome LIKE '%(drone ship)' AND PAYLOAD\_MASS\_\_KG\_ > 4000 AND PAYLOAD\_MASS\_\_KG\_
   < 6000;</li>
- List the total number of successful and failure mission outcomes. **SELECT CASE WHEN Mission\_Outcome LIKE '%success%' THEN 'Successful' ELSE 'Failures' END AS OutcomeCategory, COUNT(\*) AS TotalCount FROM SPACEXTBL GROUP BY OutcomeCategory;**
- List the names of the booster\_versions which have carried the maximum payload mass. **SELECT Booster\_Version FROM SPACEXTBL WHERE PAYLOAD\_MASS\_\_KG\_ = ( SELECT MAX(PAYLOAD\_MASS\_\_KG\_) FROM SPACEXTBL )**;
- 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. SELECT substr(Date, 6,2), Mission\_Outcome, Booster\_Version, Launch\_Site FROM SPACEXTBL WHERE substr(Date, 0,5)='2015'
- Pank the count of landing outcomes in descending order. SELECT Landing\_Outcome, COUNT(\*) AS OutcomeCount FROM SPACEXTBL WHERE Date
  BETWEEN '2010-06-04' AND '2017-03-20' GROUP BY Landing\_Outcome ORDER BY OutcomeCount DESC;
  12

Github: https://github.com/NikkideVries/IBM-SpaceX-Capstone/blob/main/4EDA\_with\_SQL.ipynb

### Build an Interactive Map with Folium

 An interactive map was created in folium
 Markers: to allow for different Exploratory Analyses and to answer questions visually.

- - Indicate launch site points to assist with identification.
- Circles:
  - Indicated highlighted regions around launch sites to help locate areas visually.
- Marker Clusters:
  - Specific launch sites in circles to help with identification
- Lines:
  - Show the distance between the coordinates

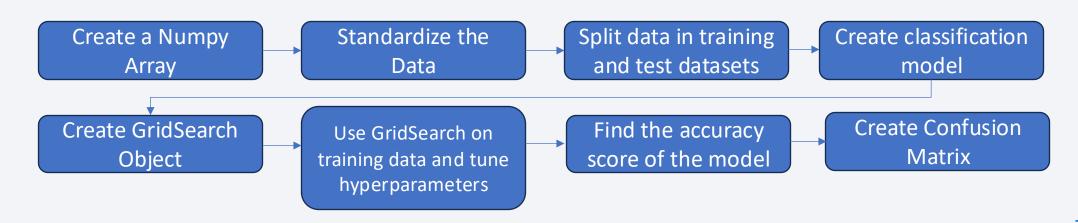
Github: https://github.com/NikkideVries/IBM-SpaceX-Capstone/blob/main/6Interactive Visuals with F olium.ipynb

# Build a Dashboard with Plotly Dash

- Using a Python interactive dashboard, Plotly Dash, an interactive dashboard was created to allow stakeholders to investigate the impacts of the launch site, payload mass, and booster type on the outcomes of the Falcon9 launch.
- Graphs that were created and that could be changed through the use of a dropdown menu were:
- A Pie chat:
  - To show which sites launched the most, and by site, what was the percentage of successful launches.
- Scatter plot:
  - To show the relationship between play load ranges and selected sites on successful launches.

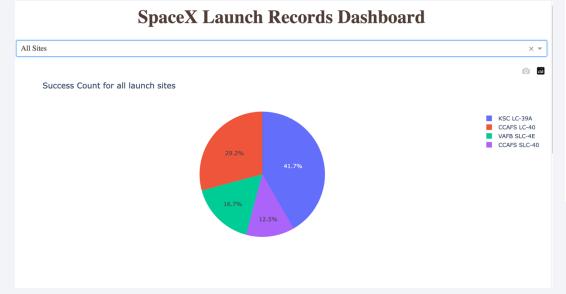
# Predictive Analysis (Classification)

- Following the flow chart, four classification four models were used: Support Vector Machine, Decision Tree, Logistic Regression, and K-nearest neighbors.
- Accuracy was used to determine the best model. All models performed the same on the test data set with an accuracy of 83.3%. However, the model with the best training accuracy was the Decision Tree model with an accuracy of 87.5%.
- Github: <a href="https://github.com/NikkideVries/IBM-SpaceX-Capstone/blob/main/8Machine\_Learning\_Prediction.ipynb">https://github.com/NikkideVries/IBM-SpaceX-Capstone/blob/main/8Machine\_Learning\_Prediction.ipynb</a>

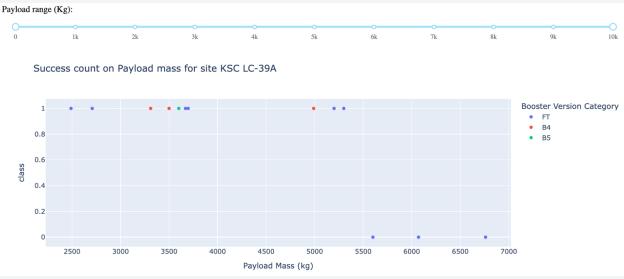


### Results

- Exploratory data analysis results
  - Launch success rate increases over time, with a large increase after 2013
  - The KSC site has the highest percentage of successful landings
  - Heavier payloads tend to lead to more landing failures

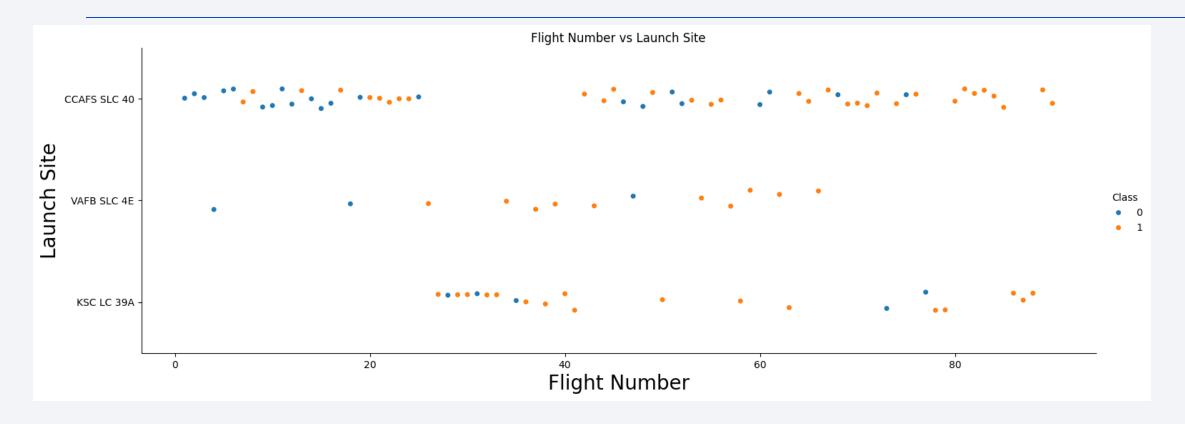


- Predictive analysis results
  - The decision tree model has the best training accuracy at 87.5%
  - All models perform the same on the test data with an accuracy of 83.3%.



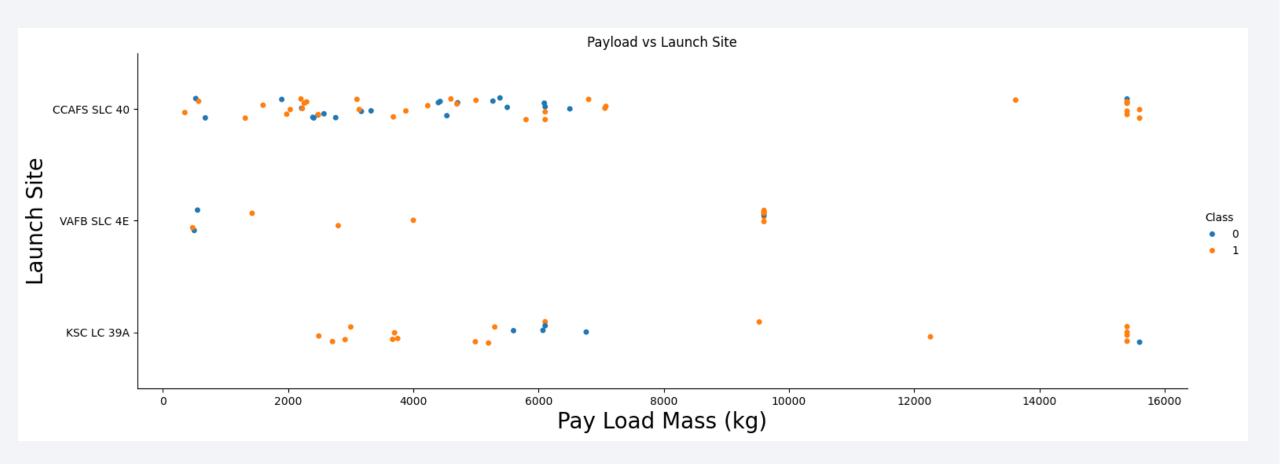


# Flight Number vs. Launch Site



- Blue class indicates landing failure while Orange class indicates landing success.
- As the flight number increases across all launch site, there are more successful landings.

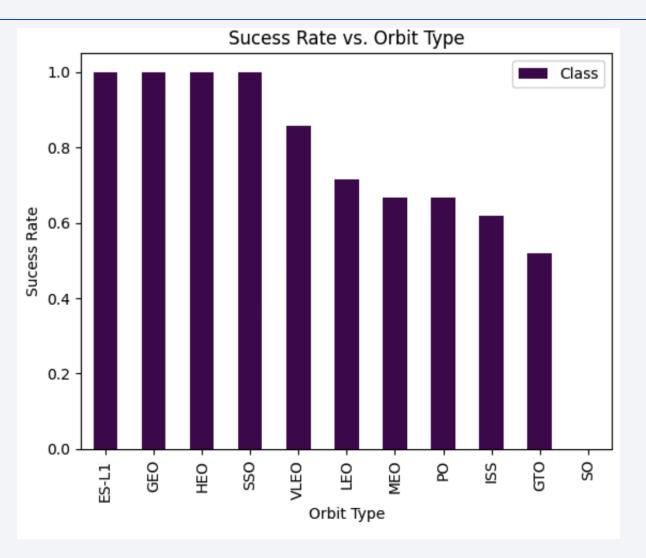
### Payload vs. Launch Site



• The higher the payload mass for sites KSC and CCA have higher successful landings.

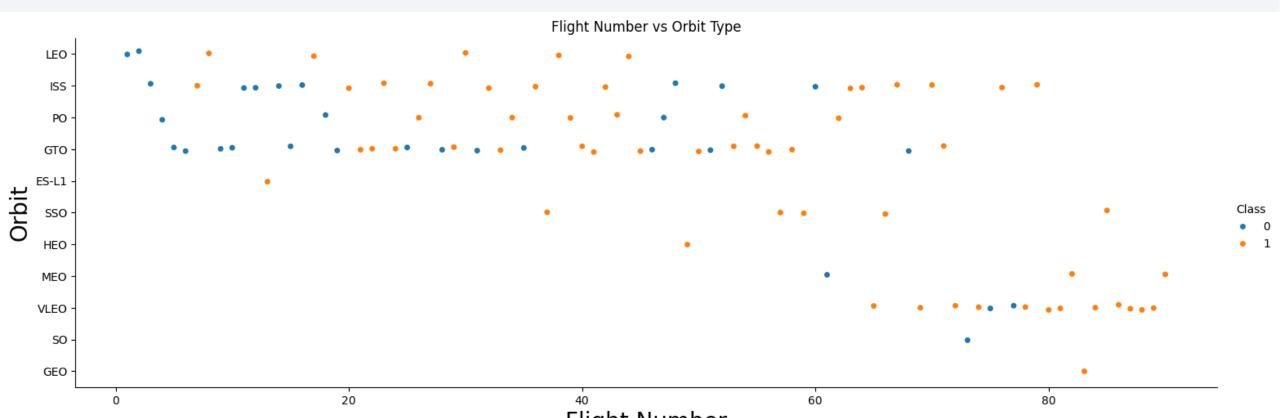
# Success Rate vs. Orbit Type

- ES-L1, GEO, HEO, and SSO orbit types have the highest successful rates.
- GTO has the lowest success rate.



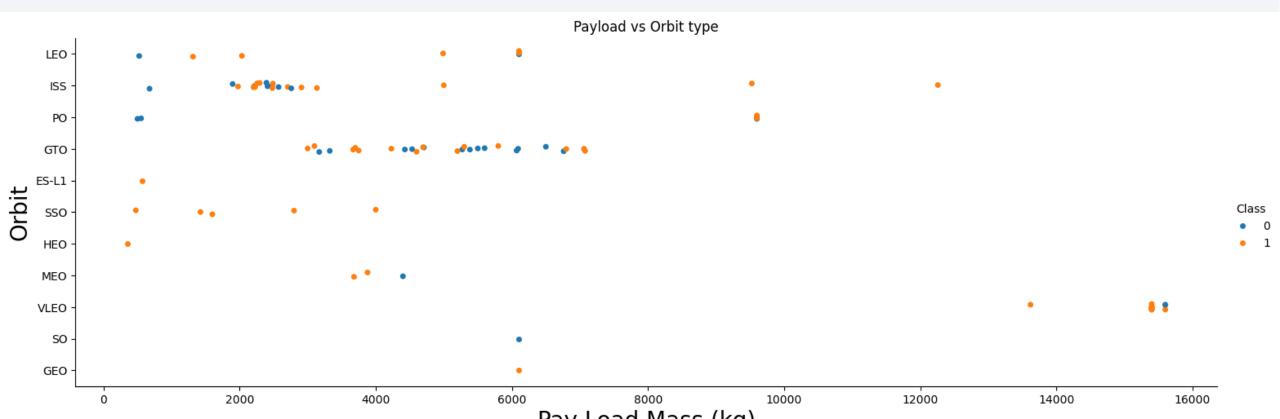
# Flight Number vs. Orbit Type

- As the number of flights increases, the number of successful landings increases.
- Blue indicates a failed landing, while Orange indicates a successful landing.



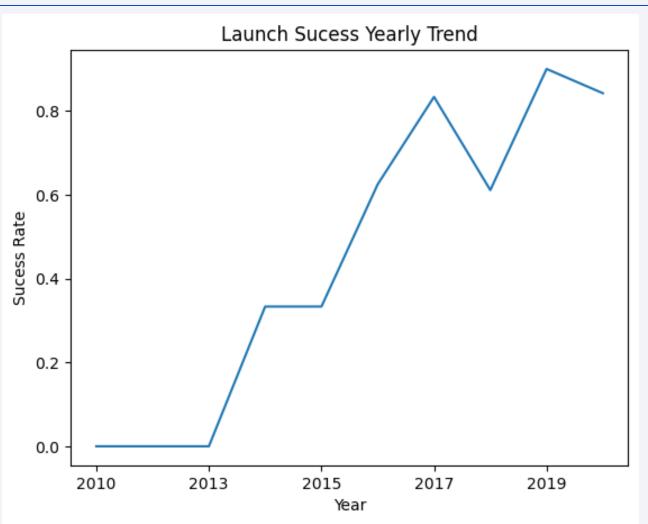
# Payload vs. Orbit Type

- There is no noticeable trend for payload vs orbit type. I would further explore this using clusters.
- Blue indicates landing failure, while Orange indicates landing success.



# Launch Success Yearly Trend

 After 2013, the success rate has steadily been increasing until 2018, where there was a dip.



### All Launch Site Names

- Display of the unique launch sites
  - SELECT DISTINCT Launch\_site FROM SPACEXTBL;
- SELECT DISTINCT will select the unique launch site names from the launch site column

#### Launch\_Site

CCAFS LC-40

VAFB SLC-4E

KSC LC-39A

CCAFS SLC-40

# Launch Site Names Begin with 'CCA'

- Display 5 records where the launch site begins with 'CCA'
  - SELECT \* FROM SPACEXTBL WHERE Launch\_Site LIKE 'CCA%' LIMIT 5;
- The command filters WHERE launch sites have names LIKE CCA, and limits it to 5 records

Date	Time (UTC)	Booster_Vers ion	Launch_Site	Payload	PAYLOAD_M ASSKG_	Orbit	Customer	Mission_Out come	Landing_Out come
2010-06-04	18:45:00	F9 v1.0 B0003	CCAFS LC-40	Dragon Spacecraft Qualification Unit	0	LEO	SpaceX	Success	Failure (parachute)
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)
2012-05-22	7:44:00	F9 v1.0 B0005	CCAFS LC-40	Dragon demo flight C2	525	LEO (ISS)	NASA (COTS)	Success	No attempt
2012-10-08	0:35:00	F9 v1.0 B0006	CCAFS LC-40	SpaceX CRS-1	500	LEO (ISS)	NASA (CRS)	Success	No attempt
2013-03-01	15:10:00	F9 v1.0 B0007	CCAFS LC-40	SpaceX CRS-2	677	LEO (ISS)	NASA (CRS)	Success	No attempt

# **Total Payload Mass**

- The total payload mass carried by boosters launched by NASA
  - SUM(PAYLOAD\_MASS\_\_KG\_) AS TotalPayloadMass FROM SPACEXTBL WHERE Customer = 'NASA (CRS)'
- SUM, sums the payload values from the customers from NASA.

**TotalPayloadMass** 

45596

### Average Payload Mass by F9 v1.1

- Calculate the average payload mass carried by booster version F9 v1.1
  - SELECT AVG(PAYLOAD\_MASS\_\_KG\_) FROM SPACEXTBL WHERE Booster\_Version = 'F9 v1.1'
- The method involves AVG(averaging) the payload mass from the specific F9 v1.1 Booster

AVG(PAYLOAD\_MASS\_\_KG\_)
2928.4

# First Successful Ground Landing Date

- Find the dates of the first successful landing outcome on ground pad
  - SELECT MIN(Date) AS EarliestDate FROM SPACEXTBL WHERE Landing\_Outcome LIKE '%(ground pad)' AND Mission\_Outcome = 'Success';
- The MIN command finds the most recent date from the landing outcomes that meet the ground pad a success parameters

**EarliestDate** 

2015-12-22

#### Successful Drone Ship Landing with Payload between 4000 and 6000

- List the names of boosters which have successfully landed on drone ship and had payload mass greater than 4000 but less than 6000
  - SELECT Booster\_Version FROM SPACEXTBL WHERE
     Landing\_Outcome LIKE '%(drone ship)' AND
     PAYLOAD\_MASS\_\_KG\_ > 4000 AND PAYLOAD\_MASS\_\_KG\_ <
     6000;</li>
- The code asks that the answer is larger than 4000 and smaller than 6000. An alternative way to solve this to use the BETWEEN command.

# F9 FT B1020 F9 FT B1022

F9 FT B1026

F9 FT B1021.2

F9 FT B1031.2

#### Total Number of Successful and Failure Mission Outcomes

Calculate the total number of successful and failure mission outcomes:

**SELECT** 

**CASE** 

WHEN Mission\_Outcome LIKE '%success%' THEN 'Successful'

**ELSE 'Failures'** 

END AS OutcomeCategory,

COUNT(\*) AS TotalCount

FROM SPACEXTBL

GROUP BY OutcomeCategory;

• The case statement checks if the mission outcome contains the string success, then labels it as outcome category (success, failure) and counts the total number of records in each category.

TotalCount	OutcomeCategory
1	Failures
100	Successful

# **Boosters Carried Maximum Payload**

• List the names of the booster which have carried the maximum payload mass:

```
SELECT Booster_Version

FROM SPACEXTBL

WHERE PAYLOAD_MASS__KG_ = (

SELECT MAX(PAYLOAD_MASS__KG_)

FROM SPACEXTBL

);
```

• The subquery finds the maximum payload mass from the table and then becomes filtered to find the booster version that has that mass.

#### **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

- List the failed landing\_outcomes in drone ship, their booster versions, and launch site names for in year 2015
  - SELECT Booster\_Version, Launch\_Site, Landing\_Outcome FROM SPACEXTBL WHERE Landing\_Outcome LIKE 'Failure (drone ship)' AND Date BETWEEN '2015-01-01' AND '2015-12-31';
- Filter through the data to find landing outcomes like failure and between year of 2015.

Booster_Version	Launch_Site	Landing_Outcome
F9 v1.1 B1012	CCAFS LC-40	Failure (drone ship)
F9 v1.1 B1015	CCAFS LC-40	Failure (drone ship)

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

 Rank the count of landing outcomes (such as Failure (drone ship) or Success (ground pad)) between the date 2010-06-04 and 2017-03-20, in descending order

SELECT Landing\_Outcome, COUNT(\*) AS OutcomeCount

FROM SPACEXTBL

WHERE Date BETWEEN '2010-06-04' AND '2017-03-20'

GROUP BY Landing\_Outcome

ORDER BY OutcomeCount DESC;

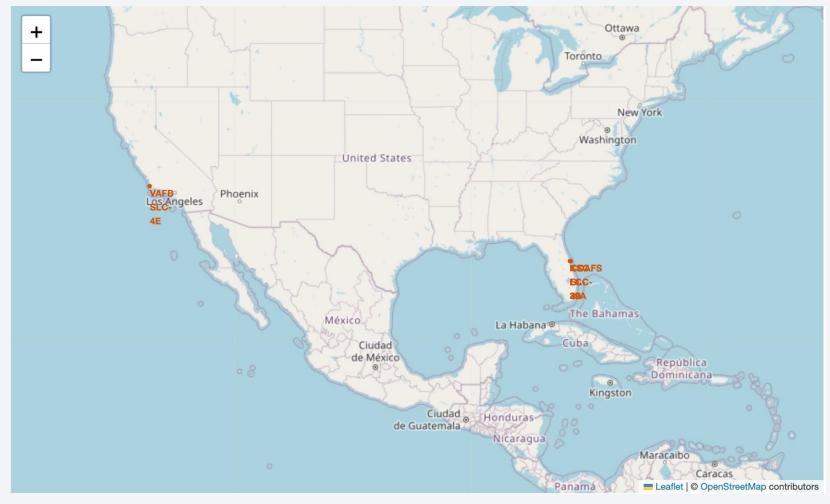
• The where statement selects for the dates, group by groups by landing outcome, and order by put them in descending order.

Landing_Outcome	OutcomeCount
No attempt	10
Success (drone ship)	5
Failure (drone ship)	5
Success (ground pad)	3
Controlled (ocean)	3
Uncontrolled (ocean)	2
Failure (parachute)	2
Precluded (drone ship)	1



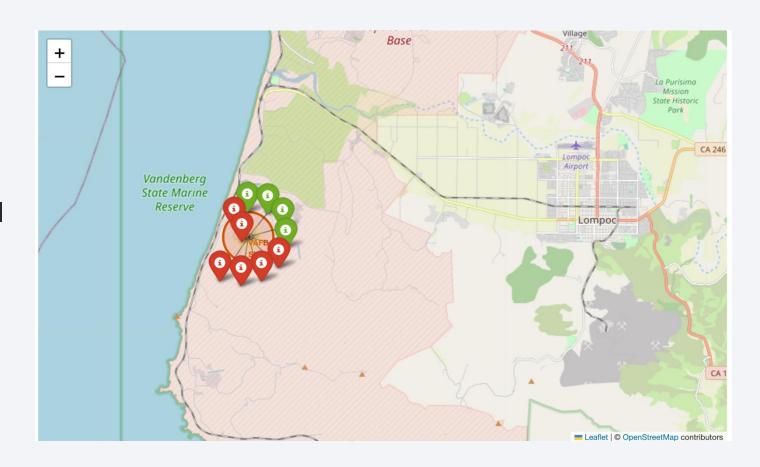
### All Launch Sites

- Vanderberg Space
   Force Base (VAFB) is
   located in California.
- Cape Canaveral (CCA)
   and Kennedy Space
   Center (KSC) are both
   located in Florida.



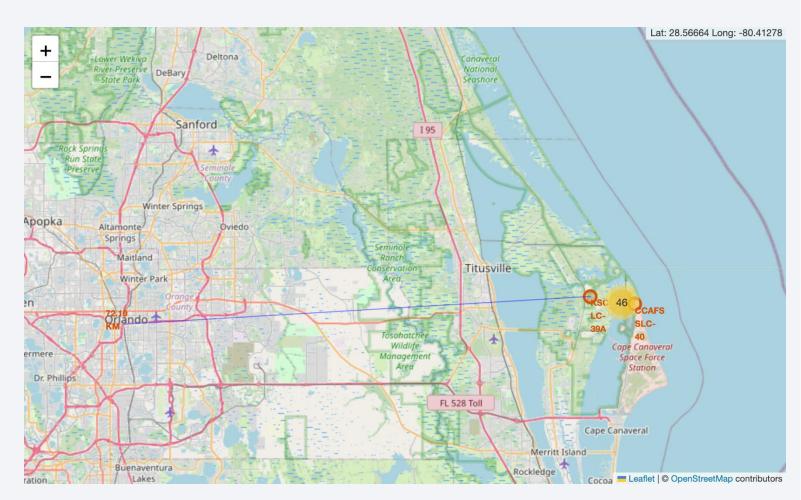
# Successful Landings Within a Launch Site

- A close-up view of the landing outcomes at Vandenberg Space Force Base.
- Red indicates an failed landing and green indicates a successful landing.



### Distance Site Proximities

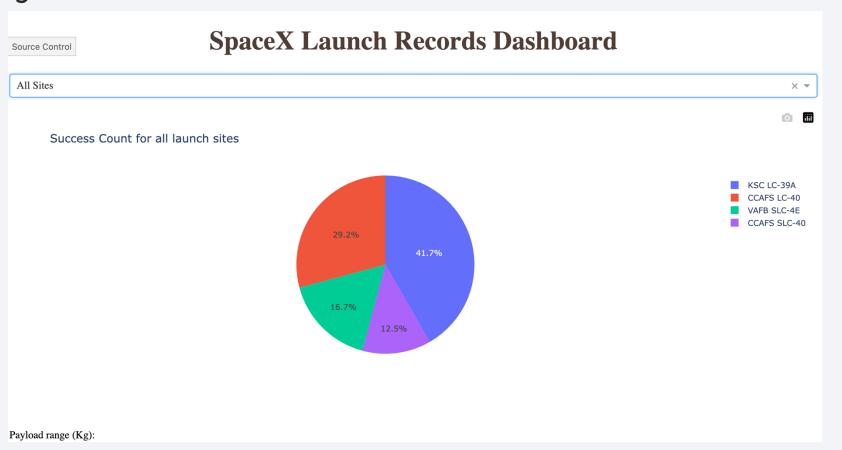
- Proximity of Kennedy Space Center (KSC) to Orlando.
- The distance is 72.10 Km.





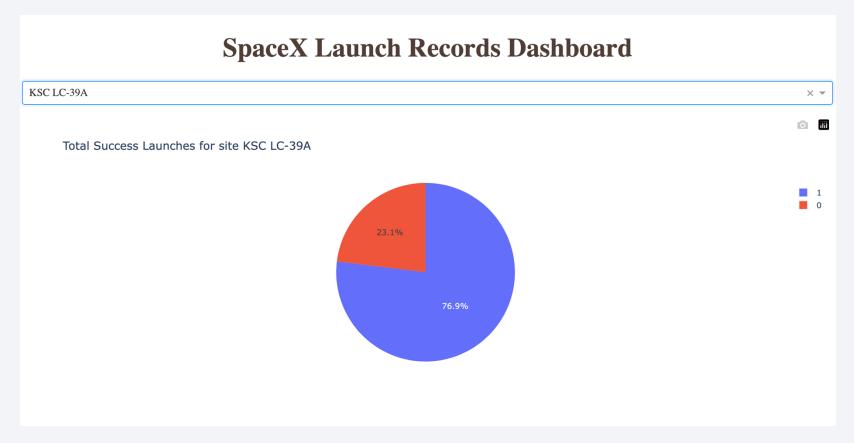
### Launch Success for all Launch Locations

 Kennedy Space Center (KPC) has the highest percentage of successful landings.



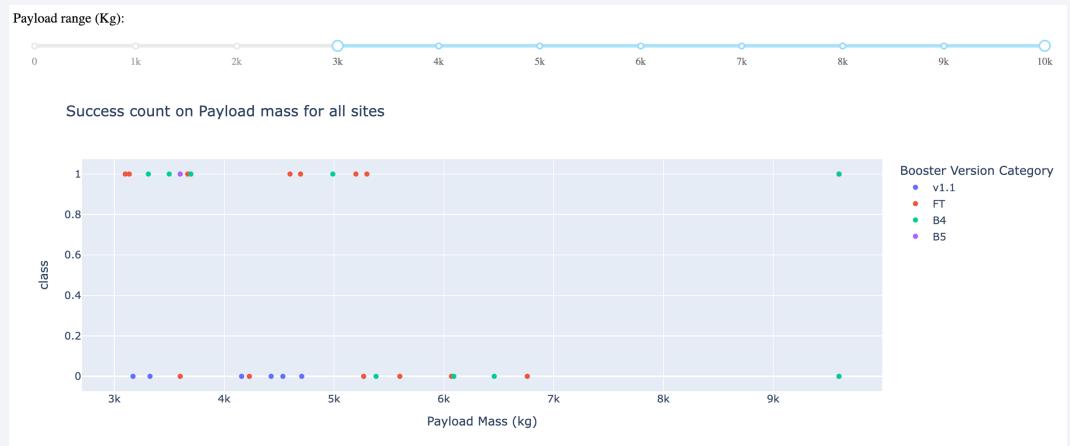
# Success Rate for the Highest Success Location KSC

 The Kennedy Space Center Launch Pad has a successful landing for Falcon 9 mission at 76.9 percent.



### Influence of Payload Mass on Success Rates

- The FT booster has the most overall Success.
- Lower payload masses seem to lead to a higher chance of successful landings.

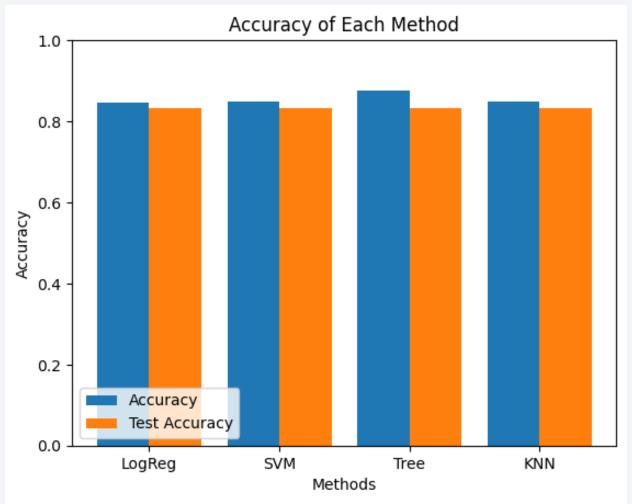




### Classification Accuracy

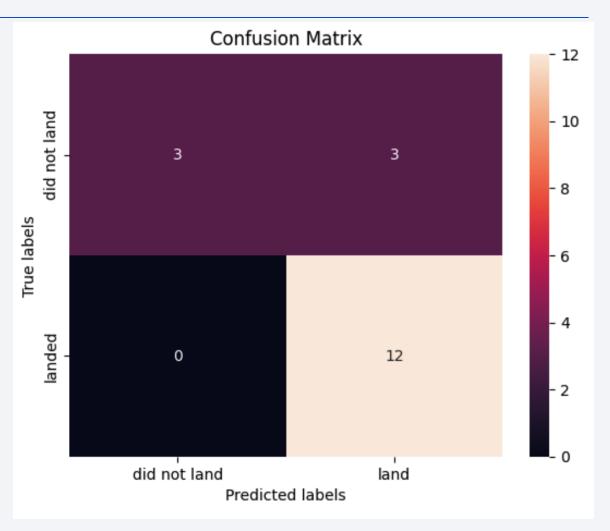
- The best training accuracy mode is the Decision tree, with a score of 87.5%.
- All models, however, had the same Test Accuracy at 83.3%.

Model	Accuracy	TestAccuracy
LogReg	0.84643	0.8333333333333334
SVM	0.84821	0.8333333333333334
Tree	0.875	0.8333333333333334
KNN	0.84821	0.83333333333333334



### **Confusion Matrix**

- True Positive (TP): 12 cases where the decision tree correctly predicted: "landed."
- True Negative (TN): 3 cases where the decision tree correctly predicted: "did not land".
- False Positive (FP): 3 cases where the decision tree incorrectly predicted "landed" when it was "did not land".
- False Negative (FN): O cases where the decision tree incorrectly predicted "did not land" when it was "landed".
- The accuracy of this model was 83.3%



### **Conclusions**

- Successful landing rates for Falcon 9 Rockets have increased over time. This could be due to better technological advancements.
- The most successful launch site is Kennedy Space Center (SPC), with a success rate of 76.9%.
- The best model is the Decision Tree, with a training accuracy of 87.5% and a test accuracy of 83.3%.

# **Appendix**

- Coursera: <a href="https://www.coursera.org/professional-certificates/ibm-data-science">https://www.coursera.org/professional-certificates/ibm-data-science</a>
- Repository: <a href="https://github.com/NikkideVries/IBM-SpaceX-Capstone/tree/main">https://github.com/NikkideVries/IBM-SpaceX-Capstone/tree/main</a>
- GitHub links:
  - Collection API: <a href="https://github.com/NikkideVries/IBM-SpaceX-Capstone/blob/main/1Data Collection API.ipvnb">https://github.com/NikkideVries/IBM-SpaceX-Capstone/blob/main/1Data Collection API.ipvnb</a>
  - Collection Web Scrapping: https://github.com/NikkideVries/IBM-SpaceX-Capstone/blob/main/2Data Collection with Web Scrapping.ipynb
  - Data Wrangling: <a href="https://github.com/NikkideVries/IBM-SpaceX-Capstone/blob/main/3Data Wrangling.ipynb">https://github.com/NikkideVries/IBM-SpaceX-Capstone/blob/main/3Data Wrangling.ipynb</a>
  - EDA with SQL: https://github.com/NikkideVries/IBM-SpaceX-Capstone/blob/main/4EDA with SQL.ipvnb
  - EDA With Visualizations: https://github.com/NikkideVries/IBM-SpaceX-Capstone/blob/main/5EDA with Visualizations.ipynb
  - Visuals with Folium: https://github.com/NikkideVries/IBM-SpaceX-Capstone/blob/main/6Interactive\_Visuals\_with\_Folium.ipynb
  - Dashboard: https://github.com/NikkideVries/IBM-SpaceX-Capstone/blob/main/7Interactive\_Dashboard\_with\_Ploty.py
  - Machine Learning Predictions: https://github.com/NikkideVries/IBM-SpaceX-Capstone/blob/main/8Machine\_Learning\_Prediction.ipynb

