# Winning Space Race with Data Science
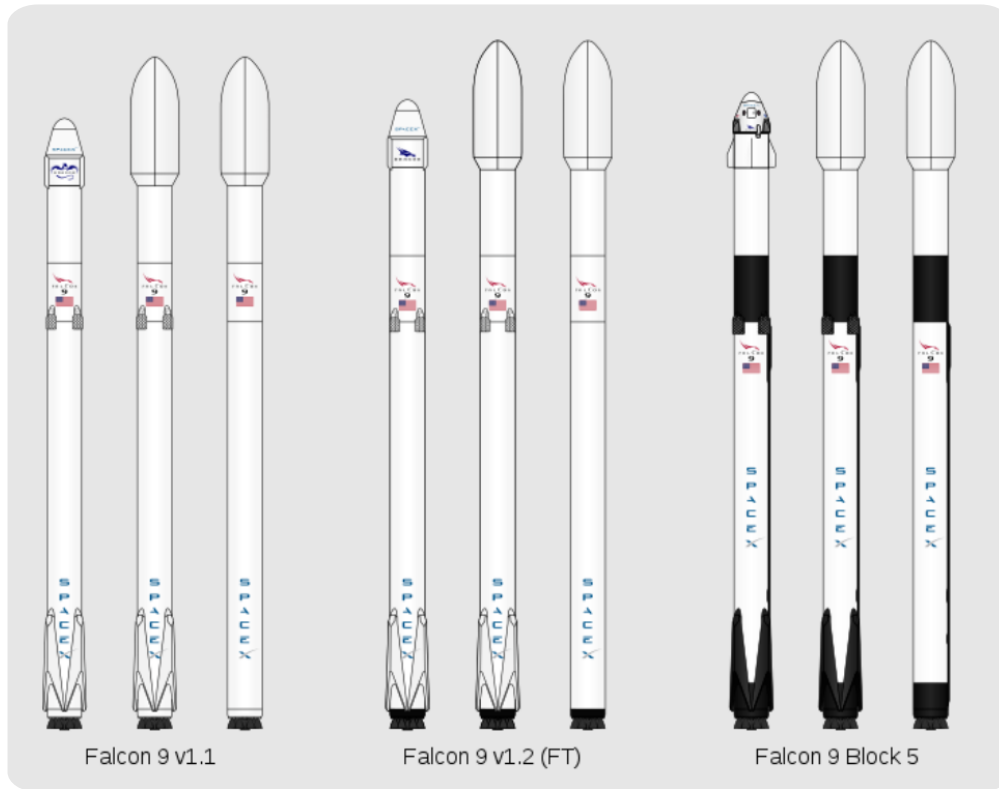
Dr Sam Thompson
October 2023

# Outline

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

# Executive Summary

- Data collected from SpaceX API and web-scraped from Wikipedia

- Data wrangled and encoded for ML

- EDA performed with visualization and SQL

  - Increase in success from 2017 to 2020, best success in ES-L1, GEO, HEO, and SSO orbit types

- Folium

  - Locations of most successful sites and proximities to landmarks

- Dashboard

  - Success rates of each launch site, most successful being CCAFS SLC-40

- ML

  - SVM and KNN produce best predictive result of all algorithms, ~84% accuracte

# Introduction



- Working for Space Y, the leading competitor to Space X, we need to understand how our rivals successes can benefit us and lead to greater market share in the reusable rocket race

- Space X's Falcon 9 rocket, their flagship vessel, has a reusable first stage. The first stage serves as the main thruster out of the atmosphere, and is extremely costly to build from scratch for each launch

- In this report we'll look at the cost of each launch, and the probability of being able to reuse the first stage of our rocket in a later trial. To do this, we'll use existing open-access data from Space X.

Section 1

# Methodology

# Methodology



**Data collection**

The data was collected directly from Space X using the Space X REST API, api.spacexdata.com/v4/. A comparison was also made using webscraping with BeautifulSoup

**Perform data wrangling**

Data was cleaned and processed using Pandas and Numpy

**Perform exploratory data analysis (EDA) using visualization and SQL**

**Perform interactive visual analytics using Folium and Plotly Dash**

**Perform predictive analysis**

Classification models were used to build our predictive capability.

# Data Collection – SpaceX API

- Data pulled from SpaceX REST API



**1) requests.get(url*)**

**2) pd.json_normalize**

Retrieve the JSON file via the API and convert to pandas dataframe

**3) Create subset of dataframe containing only 15 relevant columns of**

**4) Use specific requests.get() functions to fill columns**

**5) Isolate only rows for Falcon 9 boosters**

**6) Convert nulls in Payload to mean for column**

7

# Data Collection - Scraping

- Comparison made with webscraping from Wikipedia table. BeautifulSoup and Pandas used to make data useable

1) requests.get(url*)

2) BeautifulSoup(url.content, htmlparser)

Retrieve the data file via the API and convert to BeautifulSoup html object

3) Extract all table objects from html data BeautifulSoup.bs.findall(table)

4) For loop to extract all column names ,<th>, from html data and create dictionaries with column names

5) Loop through html data to extract values for columns and add to the dictionaries

Convert filled dictionary to a pd.DataFrame object and save as a CSV

GITHUB: IBMDataScienceCapstone/jupyter-labs-webscraping.ipynb at main · SpongeSJT/IBMDataScienceCapstone (github.com)

# Data Wrangling

- The data were explored, using value counts and date type investigations on LaunchSite, Orbit, and Outcome columns

- A dictionary was created for the Outcome column, to isolate each type of outcome

- The column converted the types of landing to either 1 for success (True), or 0 for failure (False)

- Feature engineering was also performed on the columns "Orbit", "LaunchSite", "LandingPad", "Serial", with "one hot encoding" used to prepare the data for Machine Learning

```
[23]: for i in df['Outcome']:
          if i in set(bad_outcomes):
              landing_class.append(0)
          else:
              landing_class.append(1)
```

```
[24]: df['Class']=landing_class
      df[['Class']].head(8)
```

[24]:
| | Class |
|---|---|
| 0 | 0 |
| 1 | 0 |
| 2 | 0 |
| 3 | 0 |
| 4 | 0 |
| 5 | 0 |
| 6 | 1 |
| 7 | 1 |

```
[10]: # landing_outcomes = values on Outcome column
      landing_outcomes = df['Outcome'].value_counts()
      landing_outcomes
```

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

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

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

We create a set of outcomes where the second stage did not land successfully:

```
[12]: bad_outcomes=set(landing_outcomes.keys()[[1,3,5,6,7]])
```

# EDA with Data Visualization

- ## Charts plotted

  - Flight number vs Launch Site, Scatter plot. Allowed us to see trend of how many launches were performed per site and when they were done in the sequence of launches

  - Payload Mass vs Launch Site, Scatter plot. Allowed visualization of trend between payload mass and launch site

  - Orbit Type vs Success Rate, Bar plot. Allowed visualization of success rate of launches vs their desired orbit

  - Flight Number vs Orbit type, Scatter plot. Allowed visualization of changes in orbit type during the sequence of launches

  - Payload vs Orbit type, Scatter plot. Allowed visualization of different payload masses and their desired orbit

  - Yearly Success rate change, Line plot. Showing how success of launches changed over time

GITHUB: IBMDataScienceCapstone/vizeda.ipynb at main · SpongeSJT/IBMDataScienceCapstone (github.com)

# EDA with SQL

- ## SQL Queries:

  - Display unique launch site names

  - Display sites containing string 'CCA'

  - Display sum of payload masses launched in the table

  - Display average payload mass in the table

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

GITHUB: [IBMDataScienceCapstone/sqleda.ipynb at main · SpongeSJT/IBMDataScienceCapstone (github.com)](github.com)

# Build an Interactive Map with Folium

- Explain why you added those objects

- Add the GitHub URL of your completed interactive map with Folium map, as an external reference and peer-review purpose

- Created a map of the SpaceX launch sites

- Added map markers to show their locations

- Added circles to show the number of launches per site

- Added lines to show the nearest city, the nearest railroad, and the nearest highway

# Build a Dashboard with Plotly Dash

- A Plotly Dashboard was created with pie charts and scatter plots

- The pie charts were plotted to show the number of successful landings across the different launch sites

- The pie chart could show all of the launch sites and their share of the total successes, or the success rate of a single site

- The scatter plot was made to show the relationship between payload mass, booster version, and whether or not the landing was successful

- The scatter plot was controllable by a payload mass slider to isolate regions of interest in the data space

GITHUB: [IBMDataScienceCapstone/dash.txt at main · SpongeSJT/IBMDataScienceCapstone (github.com)](IBMDataScienceCapstone/dash.txt at main · SpongeSJT/IBMDataScienceCapstone (github.com))

# Predictive Analysis (Classification)

- Four machine learning algorithms were explored
  - K Nearest Neighbors
  - Logistic Regression
  - Support Vector Machine (SVM)
  - Decision Tree
- Data split into X and Y, and test train split. X was transformed with a standard scaler prior to test train splitting
- Model was trained on a training set of the data comprising of 80% of the original dataset
- GridSearchCV was used to tune the model hyperparameters
- Accuracy on training set was returned, then accuracy on the testing set. A confusion matrix was then plotted
- K Nearest Neighbors and SVM returned the highest accuracy

Data split into X and Y, X transformed with standardscalar → Model fit with training set and GridSearchCV performed to find hyperparameters

Accuracy Score retrieved on training set → Accuracy Score retrieved on test set

Confusion matrix plotted to test overall performance of algorithm

# Insights drawn from EDA

# Flight Number vs. Launch Site

- These plots show us the number of launches per site over the span of the project

- We can see when in sequence sites were most active

- We can get a sense of success rate by class, 1= success, 0 = failure

# Payload vs. Launch Site

- This plot shows us the number of launches per site with a given payload mass

- We can see the distribution of weights launched and see any preference between site and payload mass. Highest masses are reserved for KSC and CCAFS sites.

- We can get a sense of success rate by class, 1= success, 0 = failure



Relationship between Payload Mass (kg) and Launch Site

# Success Rate vs. Orbit Type

- This plot shows us the success rate by Orbit type

- We can see that the most successful launch types are ES-L1, GEO, HEO, and SSO.

- Least successful being SO, GTO, and ISS

- No information about the distribution of sites and orbit type



Success rate of launch by Orbit type

# Flight Number vs. Orbit Type

- We can see the distribution of orbit types over the span of the data

- We see consistent launches in LEO, ISS, PO, and GTO

- GEO, SO, VLEO, MEO only occur later

- Success driven by more experience?



Relationship between Orbit type and Flight Number

# Payload vs. Orbit Type

- We see overall a preference of <7000kg Payload mass

- GTO has a narrow band of launches between 2000 and 8000 kg

- Smallest ranges are in MEO

- Largest range in ISS



Relationship between Orbit type and Payload mass (kg)

# Launch Success Yearly Trend

- We see a steady increase in success rate with a slight dip towards 2020

- This could be due to fewer launches in the latter years of the dataset



Success rate of launches over time

# All Launch Site Names

- Find the names of the unique launch sites

- We see the list of 4 launch sites

```
[10]: %sql SELECT DISTINCT "Launch_Site" FROM SPACEXTBL

    * sqlite:///my_data1.db
Done.
[10]: ,,,,,,,,,,,,,,,,,,
```

**Launch_Site**

CCAFS LC-40

VAFB SLC-4E

KSC LC-39A

CCAFS SLC-40

# Launch Site Names Begin with 'CCA'

- Find 5 records where launch sites begin with `CCA`

- We see that there are 60 launches from sites beginning with CCA

```
[11]: %sql SELECT count(*) from SPACEXTBL Where "Launch_Site" LIKE '%CCA%'

       * sqlite:///my_data1.db
```

Done.

[11]: , , , , , , , , , , ,

**count(*)**

60

# Total Payload Mass

- Calculate the total payload carried by boosters from NASA

- The total sum of the payloads is displayed here as 619967 kg

## Task 3

Display the total payload mass carried by boosters launched by NASA (CRS)

```
[12]: %sql SELECT SUM("PAYLOAD_MASS__KG_") FROM SPACEXTBL
```

 * sqlite:///my_data1.db

Done.

[12]: ,,,,,,,,,,

**SUM("PAYLOAD_MASS__KG_")**

619967

# Average Payload Mass by F9 v1.1

- Calculate the average payload mass carried by booster version F9 v1.1

- We see the average mass carried is 2928.4 kg

```
[13]: %sql SELECT AVG("PAYLOAD_MASS__KG_") FROM SPACEXTBL WHERE "Booster_Version" like 'F9 v1.1'
       * sqlite:///my_data1.db
      Done.
[13]: ,,,,,,,,,,
      AVG("PAYLOAD_MASS__KG_")
                        2928.4
```

# First Successful Ground Landing Date

- Find the dates of the first successful landing outcome on ground pad

- The earliest success in our data set was on the 22nd of December 2015

```
[14]: %sql SELECT MIN("Date") from SPACEXTBL WHERE "Landing_Outcome" like '%success%'

      * sqlite:///my_data1.db
Done.
[14]: ,,,,,,,,,,,

      MIN("Date")

      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

- We see that only F9 FT boosters have carried weights between 4000 and 6000 kg

```
[15]: %sql SELECT Booster_Version from SPACEXTBL WHERE "Landing_Outcome" like 'Success (drone ship)' and "PAYLOAD_MASS__KG_" between 4000 and 6000

 * sqlite:///my_data1.db
Done.
```

[15]: ,,,,,,,,,,,,,,,,,,,,

| Booster_Version |
| --- |
| 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

- We see here there have been 61 total successes

```
[16]: %sql SELECT count("Mission_Outcome") from SPACEXTBL WHERE "Landing_Outcome" like '%success%'

 * sqlite:///my_data1.db
Done.
```

[16]: ,,,,,,,,,,,

**count("Mission_Outcome")**

61

# Boosters Carried Maximum Payload

- List the names of the booster which have carried the maximum payload mass

- We see only the F9 B5 boosters have carried the maximum payload

```
[21]: %sql SELECT DISTINCT("Booster_Version") From SPACEXTBL Where "PAYLOAD_MASS__KG_"=(SELECT MAX("PAYLOAD_MASS__KG_") from SPACEXTBL)

 * sqlite:///my_data1.db
Done.
```

[21]:

| 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

```
[20]: %sql SELECT "Date", "Landing_Outcome", "Booster_Version", "Launch_Site" from SPACEXTBL Where "Date" like '%2015%' and "Landing_Outcome" like '%failure%'
 * sqlite:///my_data1.db
Done.
```

| Date | Landing_Outcome | Booster_Version | Launch_Site |
|------|-----------------|-----------------|-------------|
| 2015-10-01 | Failure (drone ship) | F9 v1.1 B1012 | CCAFS LC-40 |
| 2015-04-14 | Failure (drone ship) | F9 v1.1 B1015 | CCAFS LC-40 |

# 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

- The most prevalent result is no attempt, with a 3 way split on 2nd place for success (ground pad), success (drone ship), and failure (drone ship)

```
[23]: %sql SELECT "Landing_Outcome" ,COUNT(*) as COUNT_LAUNCHES FROM SPACEXTBL WHERE "Date" between '2010-06-04' and '2017-03-20' GROUP BY "Landing_Outcome" ORDER BY COUNT_LAUNCHES DESC
       * sqlite:///my_data1.db
      Done.
```

[23]:

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

Section 3

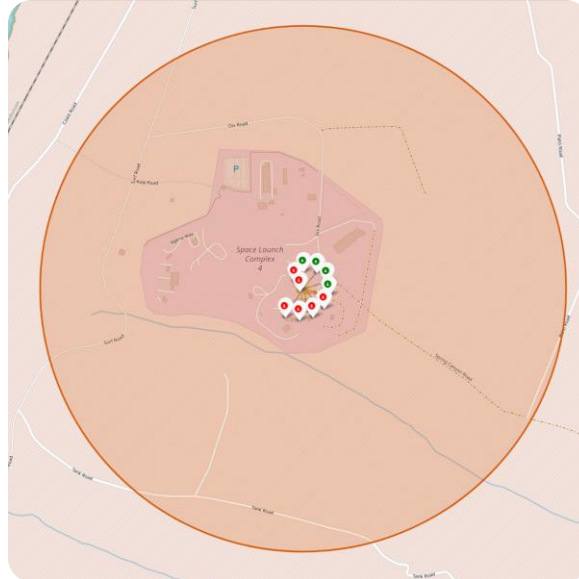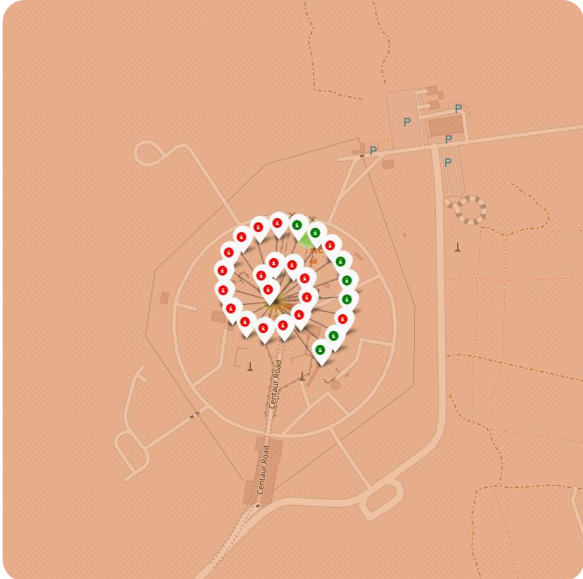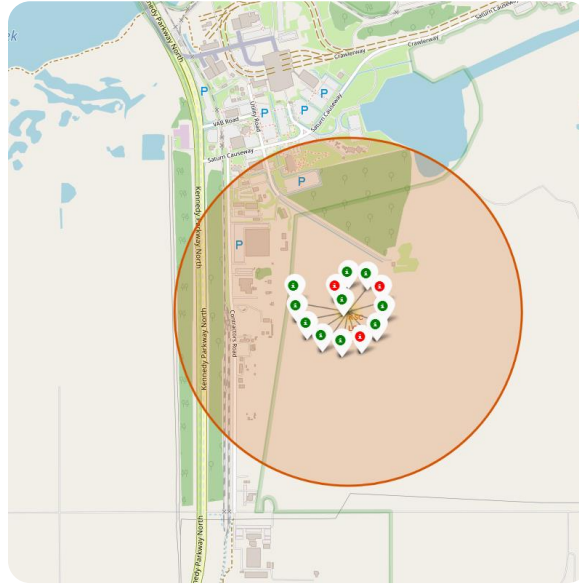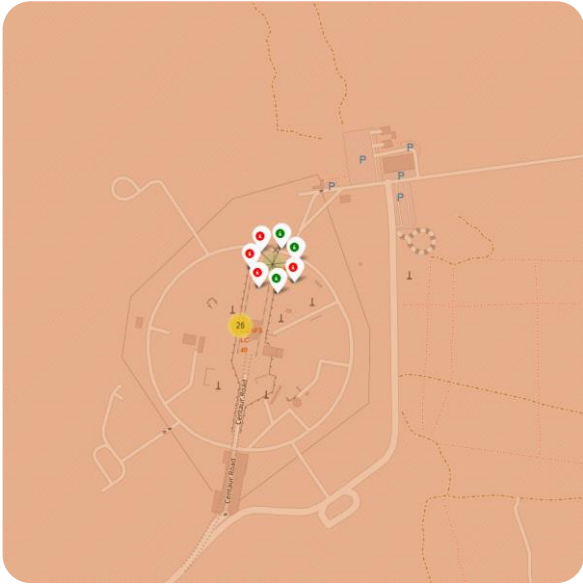# Launch Sites
# Proximities Analysis

# Launch site locations map

- We see the launch locations on the map

- They are based on the coast of Florida and California

- California was likely chosen for existing infrastructure and ease of hiring top engineering talent near the silicon valley area

- Florida likewise has a lot of existing infrastructure, and close ties for NASA. More launches were performed here than California
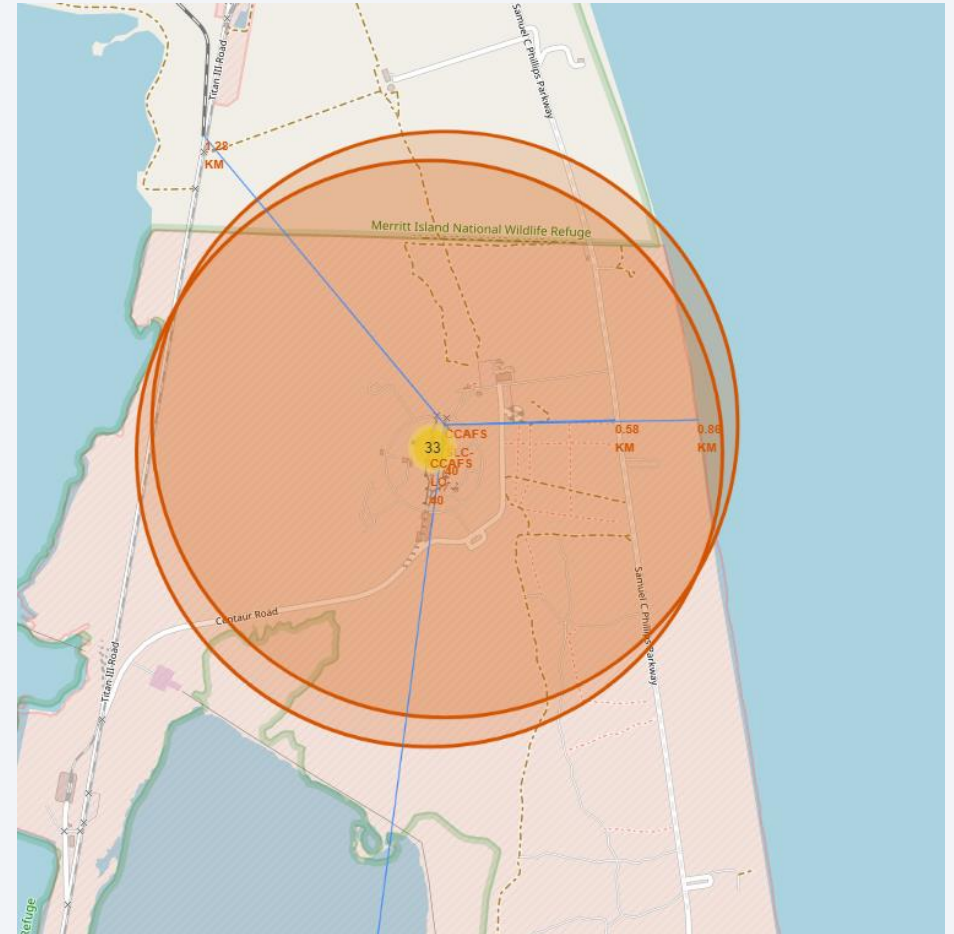
# Success and number per site

- We see the number of launches per site as a spiral of markers per site

- Red is failure, green is success

# Launch site proximity to landmarks

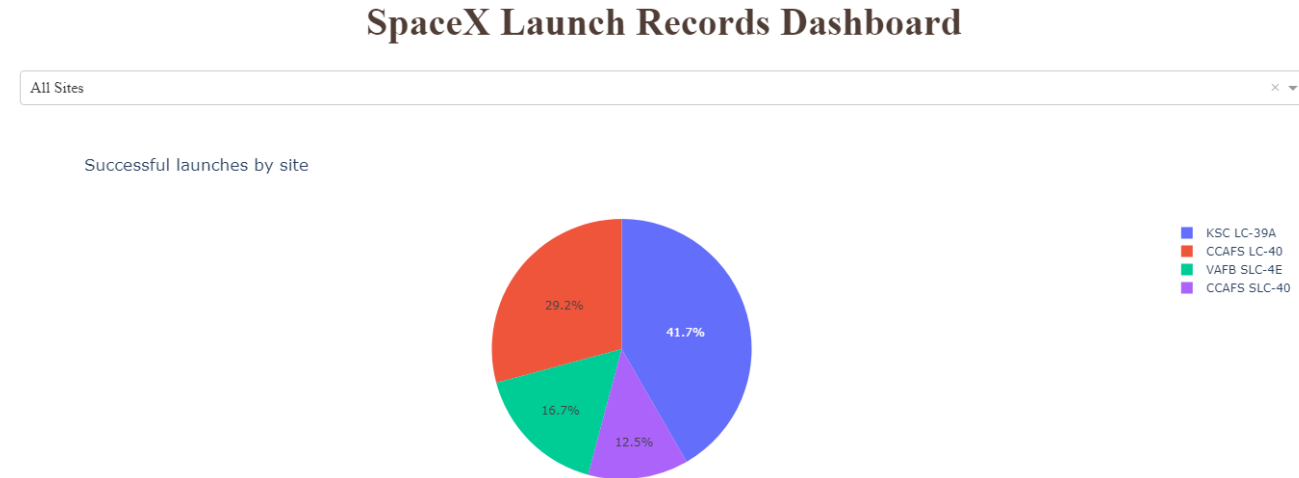- Map showing proximity of nearest railway, highway, and coastline to the launch site

Section 4

# Build a Dashboard
# with Plotly Dash

# Total Successes by Launch site
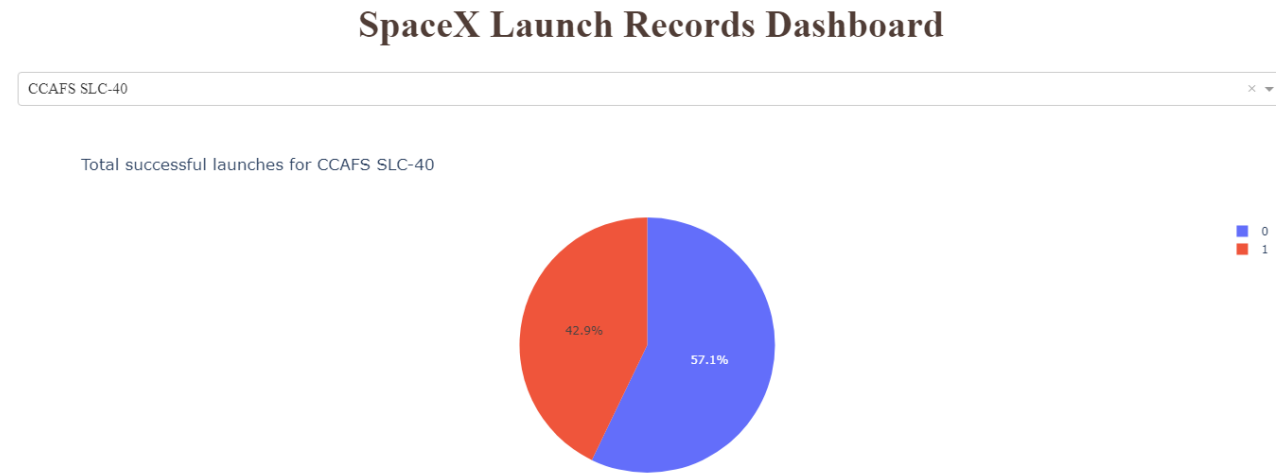
- We see here the homepage of the dashboard

- The pie chart displays the percentage of successful launches broken down by launch site

- Greatest proportion of successes attributed to KSC LC-39A

- Lowest proportion attributed to CCAFS SLC-40

**SpaceX Launch Records Dashboard**

All Sites

Successful launches by site

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

41.7%
29.2%
16.7%
12.5%

# Most successful site

- Despite making up the smallest number of successes, CCAFS SLC-40 has the highest success rate

- Fewer launches at this site but they were on the whole more successful

**SpaceX Launch Records Dashboard**

CCAFS SLC-40

Total successful launches for CCAFS SLC-40
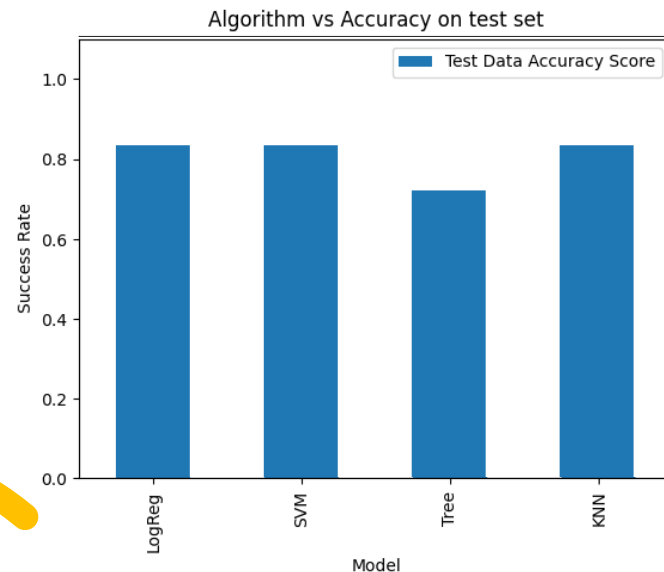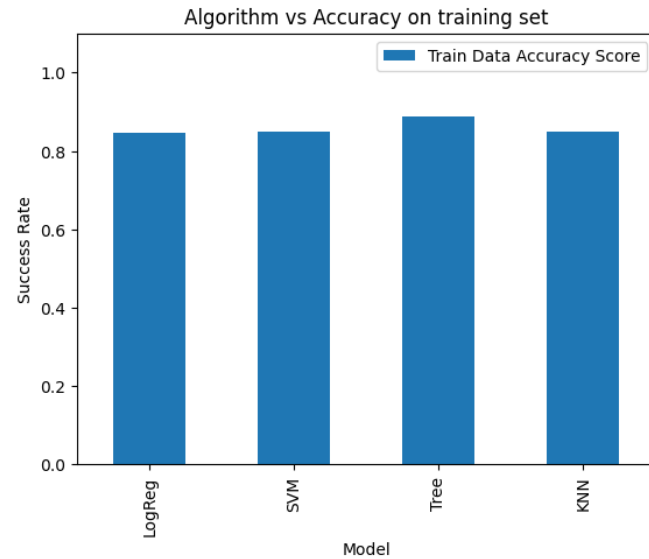
42.9%  57.1%

0
1

# Payload Mass by Site, and success rate

- Highest success rate is seen between masses of 2000 kg and 4000 kg

- The most successful boosters appear to be FT and B4

Payload range (Kg):

Section 5

# Predictive Analysis (Classification)

Algorithm vs Accuracy on training set
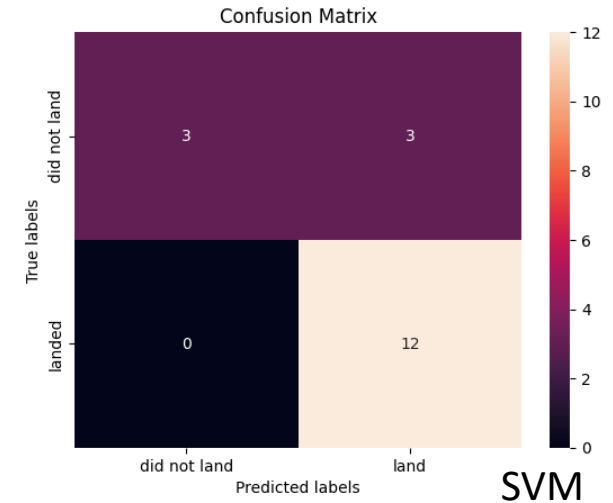


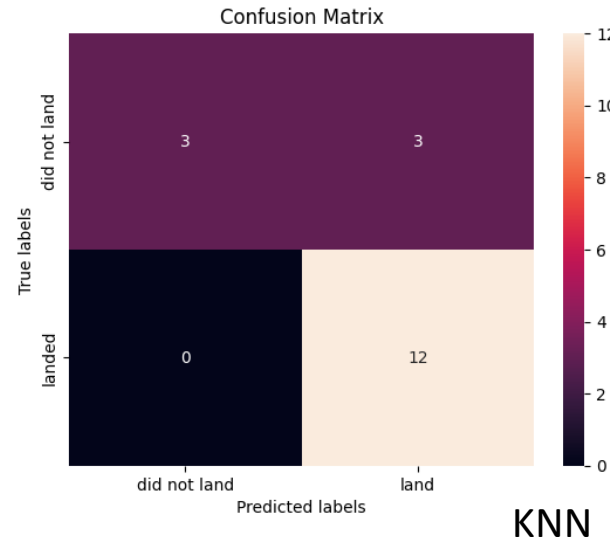Algorithm vs Accuracy on test set

# Classification Accuracy

- We see that in the training data set the decision tree algorithm out-performs all others

- However, when applied to the test set it falls flat

- The overall best performance was produced by the KNN and SVM algorithms

# Confusion Matrix

- We see the confusion matrices for the KNN and the SVM algorithms

- Both perform identically, with both producing the best accuracy scores on training and test sets

- We see no false negatives, but three false positives produced by both models



SVM



KNN

42

42

# Conclusions

- SpaceX saw a steady increase in successes until 2020

- Payload masses are most successful between 2000 and 4000 kg. we should limit our launches to these masses until high success is achieved

- The most successful site is CCAFS SLC-40, we should replicate their set-up and operations

- The classification models broadly perform to a similar level, though SVM or KNN should be used for greatest accuracy

Thank you!