**IBM Developer**
SKILLS NETWORK

# Winning Space Race with Data Science

Jason Forsythe
7/14/22

# Outline

- Executive Summary

- Introduction

- Methodology

- Results

- Conclusion

- Appendix

# Executive Summary

- Summary of methodologies

- Summary of all results

# Introduction

In this project I wanted to find out the answers to a few questions

- Can I predict the success of a SpaceX first stage landing?

- If so, what parameters could be used to give the best prediction?

- What would be the best accuracy of prediction that I could achieve?

Section 1

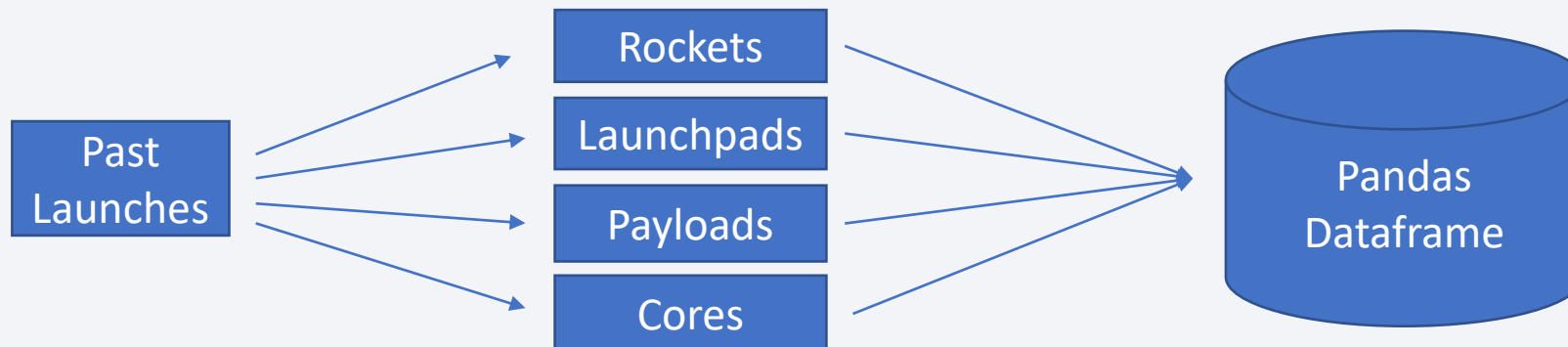# Methodology

# Methodology

Executive Summary

- Data collection methodology:

    - Data was collected via SpaceX API and Web Scraping
      https://github.com/Jason-AE/IBM_Data_Science/tree/master/DataCollectionAPI

- Perform data wrangling

    - Once Data was collected, I removed unesseary columns and replaced null values.
      https://github.com/Jason-AE/IBM_Data_Science/tree/master/DataWrangling

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

- Perform interactive visual analytics using Folium and Plotly Dash

- Perform predictive analysis using classification models

    - I tested Logistic Regression, Vector Machine, Tree, and K Nearest Neighbor models using grid search to find the best
      parameters for each. I then compared the results using test data and the Tree model performed the best.
      https://github.com/Jason-AE/IBM_Data_Science/tree/master/PredictiveAnalysis

# Data Collection – SpaceX API

- Data was collected using SpaceX's own Web API

  - https://api.spacexdata.com/v4/launches/past

  - This data was used along with different endpoints to combine data about the rockets, launchpads, payloads, and cores.

  - The Data was then stored in a Panda's Data Frame



https://github.com/Jason-AE/IBM_Data_Science/tree/master/DataCollectionAPI

# Data Collection - Scraping

- Data was also collected using Web scrapping using Wikipedia data in HTML Tables
  https://en.wikipedia.org/w/index.php?title=List_of_Falcon_9_and_Falcon_Heavy_launches&oldid=1027686922



HTML Table Example

- BeautifulSoup Library Was Used to Extract Data
- Helper functions cleaned the HTML
- Pandas Dataframe was created
- Iterations over the rows was used to fill Dataframe



Resulting Pandas Dataframe

https://github.com/Jason-AE/IBM_Data_Science/tree/master/DataCollectionAPI

# Data Wrangling

- Once data was collected it needs to be cleaned
- Check for missing values in each attribute (~40 in LandingPad)
  - Replace missing values with mean in that column or remove rows
- Convert landing outcomes to a class of either 0 failed 1 success
- Now we can check the mean of success (66.666%)
- This clean up is important for normalization and predictive analytics (to be preformed later)

https://github.com/Jason-AE/IBM_Data_Science/tree/master/DataWrangling

# EDA with Data Visualization

- In Exploratory Data Analysis I plotted the following charts:

- Success rate vs launch site: to determine if launch site affects success

- Flight number per site vs success rate: to determine if launches were more successful with later launches.

- Payload per site vs success rate: to determine if payload affects success rate at each site

- Success rate vs orbit type: to determine if orbit affects success rate

- Payload by orbit type vs success rate: to determine if certain payloads in a give orbit affect success rate

- Success rate over time (yearly): to determine if success increased over time

- https://github.com/Jason-AE/IBM_Data_Science/tree/master/EDA

# EDA with SQL

- In exploratory data analysis using SQL I performed the following queries:

- Select distinct launch site names

- Displayed 5 records from the data set where the launch site contained 'KSC'

- Displayed the total payload carried for NASA missions

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

- Showed the date of the first successful landing on a drone ship

- Listed the names of the boosters that have had successful landings on land with a mass greater than 4000 and less that 6000 kg

- Listed the total number of successful and failure mission outcomes

- Listed the booster versions that have carried the maximum payload mass

- For 2017 listed the booster versions and launch sites for successful landings on land

- Ranked the count of successful landings between 2010-06-04 and 2017-03-20 in descending order

- https://github.com/Jason-AE/IBM_Data_Science/tree/master/EDA

# Build an Interactive Map with Folium

- I built a Folium Interactive Map to show the relation of launch sites and successful landings at each in a more visual way

- I used a folium.Circle object to mark the area of the launch site with details

- I also used a folium.map.Marker to mark a launch site when zoomed out

- I made markers of green and red to represent successful or failed launches

- Lastly, I made marker clusters to group successful and failed launches based on launch sites

- https://github.com/Jason-AE/IBM_Data_Science/tree/master/FoliumMap

# Build a Dashboard with Plotly Dash

- I created an interactive dashboard with plotly with the following features:

- Pie chart showing success rate per launch site

- Scatter plot showing launch outcomes vs payload colored by booster version

- Launch site dropdown that includes all sites (Controls data for pie and scatter chart)

- Payload range slider (controls data for scatter chart only)

- https://github.com/Jason-AE/IBM_Data_Science/tree/master/Capstone_Mod3_Interactive_Dashboard

# Predictive Analysis (Classification)

- During predictive Analysis I test the following models:

- Logistic Regression

- Vector Machine

- Tree

- K Nearest Neighbor

- I used grid search to find the best parameters for each.

- I then compared the results using test data and the Tree model performed the best.
  https://github.com/Jason-AE/IBM_Data_Science/tree/master/PredictiveAnalysis

# Results

- I found that Orbit, Launch Site, Payload, and Flight Number are good predictors of success.





- Interactive analytics demo in screenshots





- Predictive analysis results



```
print("LogReg Score: ",logreg_cv.score(X_test, Y_test))
print("SVM Score: ", svm_cv.score(X_test, Y_test))
print("Tree Score: ", tree_cv.score(X_test, Y_test))
print("KNN Score: ", knn_cv.score(X_test, Y_test))
```

```
LogReg Score:  0.8333333333333334
SVM Score:  0.8333333333333334
Tree Score:  0.8888888888888888
KNN Score:  0.8333333333333334
```

15

Section 2

# Insights drawn from EDA

# Flight Number vs. Launch Site

- Scatter plot of Flight Number vs. Launch Site



- Launch Sites on the Y axis and Flight Number on the X axis, Blue dots are failed landings and Orange dots are successful landings.

- You can see that KSC LC-39A has the best record of success

- Additionally, you can see CCAFS SLC 40 has improved its success over time

# Payload vs. Launch Site

- Scatter plot of Payload vs. Launch Site



- CCAF5 SLC 40 has a better track record with heavier payload

- VAFB SLC 4E's failures happed with very light payloads

- Most success has happened with payloads from 8000 to16000 kg

# Success Rate vs. Orbit Type

- Bar chart for the success rate of each orbit type

- ESL1, GEO, HEO, and SSO have 100% success

- VLEO is close second with 90% success

- LEO is in third with 70% success

- GTO performed the worse with 50% success

- Imagine flipping a coin to see if your multi million-dollar rocket would land!

# Flight Number vs. Orbit Type

- Scatter point of Flight number vs. Orbit type



- First few launches were in LEO and ISS orbits and failed

- Launches in LEO and ISS orbits had a better success rate with time

- VLEO, SO, GEO, where only attempted later in the program

# Payload vs. Orbit Type

- Scatter point of payload vs. orbit type



- GTO orbit landing success does not seem to be tied to payload mass

- ES-L1, SSO, HEO, and MEO launches took place with less than 7000kg of payload

- VLEO launches saw the highest payload masses of the program

# Launch Success Yearly Trend

- Line chart of yearly average success rate

- The SpaceX program started in 2010

- In 2015 the success rate was less than 40%

- 2017 the success rate jumped to over 80%

- There was a slight dip in success in 2018

- 2019 rebounded with a nearly perfect landing success rate

# All Launch Site Names

- In the data set there are four unique launch sites

- CCAFS LC-40

- CCAFS SLC-40

- KSC LC-39A

- VAFB SLC-4E

- select distinct(LAUNCH_SITE) FROM SPACEXTBL

- The above query returns the distinct launch sites from the data set.

# Launch Site Names Begin with 'KSC'

| DATE | time_utc_ | booster_version | launch_site | payload | payload_mass_kg_ | orbit | customer | mission_outcome | landing_outcome |
|---|---|---|---|---|---|---|---|---|---|
| 2017-02-19 | 14:39:00 | F9 FT B1031.1 | KSC LC-39A | SpaceX CRS-10 | 2490 | LEO (ISS) | NASA (CRS) | Success | Success (ground pad) |
| 2017-03-16 | 06:00:00 | F9 FT B1030 | KSC LC-39A | EchoStar 23 | 5600 | GTO | EchoStar | Success | No attempt |
| 2017-03-30 | 22:27:00 | F9 FT B1021.2 | KSC LC-39A | SES-10 | 5300 | GTO | SES | Success | Success (drone ship) |
| 2017-05-01 | 11:15:00 | F9 FT B1032.1 | KSC LC-39A | NROL-76 | 5300 | LEO | NRO | Success | Success (ground pad) |
| 2017-05-15 | 23:21:00 | F9 FT B1034 | KSC LC-39A | Inmarsat-5 F4 | 6070 | GTO | Inmarsat | Success | No attempt |

- SELECT * FROM SPACEXTBL WHERE LAUNCH_SITE LIKE 'KSC%' limit 5

- Here are 5 records of launches from launch sites that start with KSC

# Total Payload Mass

- select sum(payload_mass__kg_) from spacextbl where customer = 'NASA (CRS)'

- The above query sums the total payload mass of all rockets where the customer was NASA

- The total payload is 45,596 kg

# Average Payload Mass by F9 v1.1

- select avg(payload_mass__kg_) from spacextbl where booster_version = 'F9 v1.1'

- The above query calculates the average payload carried by F9 v1.1 boosters

- The average payload is 2,928 kg for the F9 v1.1 boosters

# First Successful Ground Landing Date

- select min(date) from spacextbl where landing__outcome = 'Success (drone ship)'

- The above query finds the first date that SpaceX performed a successful drone ship landing

- That date was 04-08-2016

# Successful Drone Ship Landing with Payload between 4000 and 6000

```sql
select booster_version from spacextbl
where
    landing__outcome = 'Success (ground pad)'
and payload_mass__kg_ > 4000
and payload_mass__kg_ < 6000
```

| booster_version |
|---|
| F9 FT B1032.1 |
| F9 B4 B1040.1 |
| F9 B4 B1043.1 |

- The above images show the query to get the booster version that have landed on a drone ship that launched with more than 4,000kg of payload and less than 6,000kg of payload.

# Total Number of Successful and Failure Mission Outcomes

- select mission_outcome, count(*) as count from spacextbl group by mission_outcome

- The above query retrieves the counts of the different mission outcomes (this is not the same as landing outcomes of the booster)

| mission_outcome | COUNT |
|---|---|
| Failure (in flight) | 1 |
| Success | 99 |
| Success (payload status unclear) | 1 |

# Boosters Carried Maximum Payload

- select booster_version from spacextbl
  where
        payload_mass__kg_ =
  (select max(payload_mass__kg_) from spacextbl)

- The above query retrives the booster versions that have carried the maximum payload. It uses a sub query in the predicate to get the max payload.

| 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 |

# 2017 Launch Records

```
select
    monthname(date) as Month,
    landing__outcome,
    booster_version,
    launch_site
from spacextbl
where
    landing__outcome = 'Success (ground pad)'
    and year(date) = 2017
```

| MONTH | landing__outcome | booster_version | launch_site |
|---|---|---|---|
| February | Success (ground pad) | F9 FT B1031.1 | KSC LC-39A |
| May | Success (ground pad) | F9 FT B1032.1 | KSC LC-39A |
| June | Success (ground pad) | F9 FT B1035.1 | KSC LC-39A |
| August | Success (ground pad) | F9 B4 B1039.1 | KSC LC-39A |
| September | Success (ground pad) | F9 B4 B1040.1 | KSC LC-39A |
| December | Success (ground pad) | F9 FT B1035.2 | CCAFS SLC-40 |

- The above images show the query and successful landings on ground

- The columns where specifically selected to give a summary

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

```
select landing__outcome, count(*) as cnt from spacextbl
where landing__outcome like 'Success%'
group by landing__outcome
order by 2
```

| landing_outcome | cnt |
| --- | --- |
| Success (ground pad) | 9 |
| Success (drone ship) | 14 |
| Success | 38 |

- The above images show the query and the successful landings ranked by landing outcome

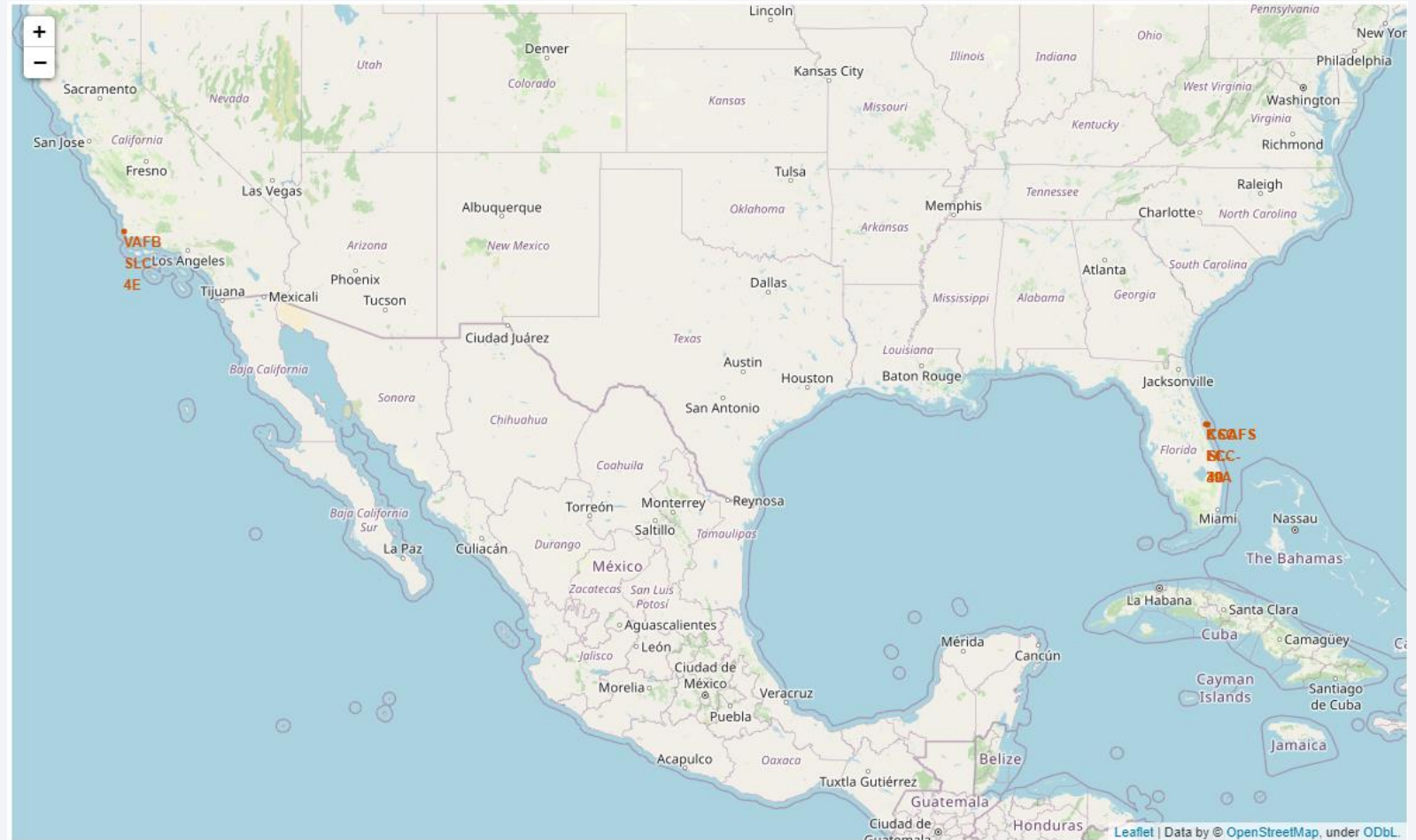- Group by must be used when an aggregate function is used such as count(*)

Section 3

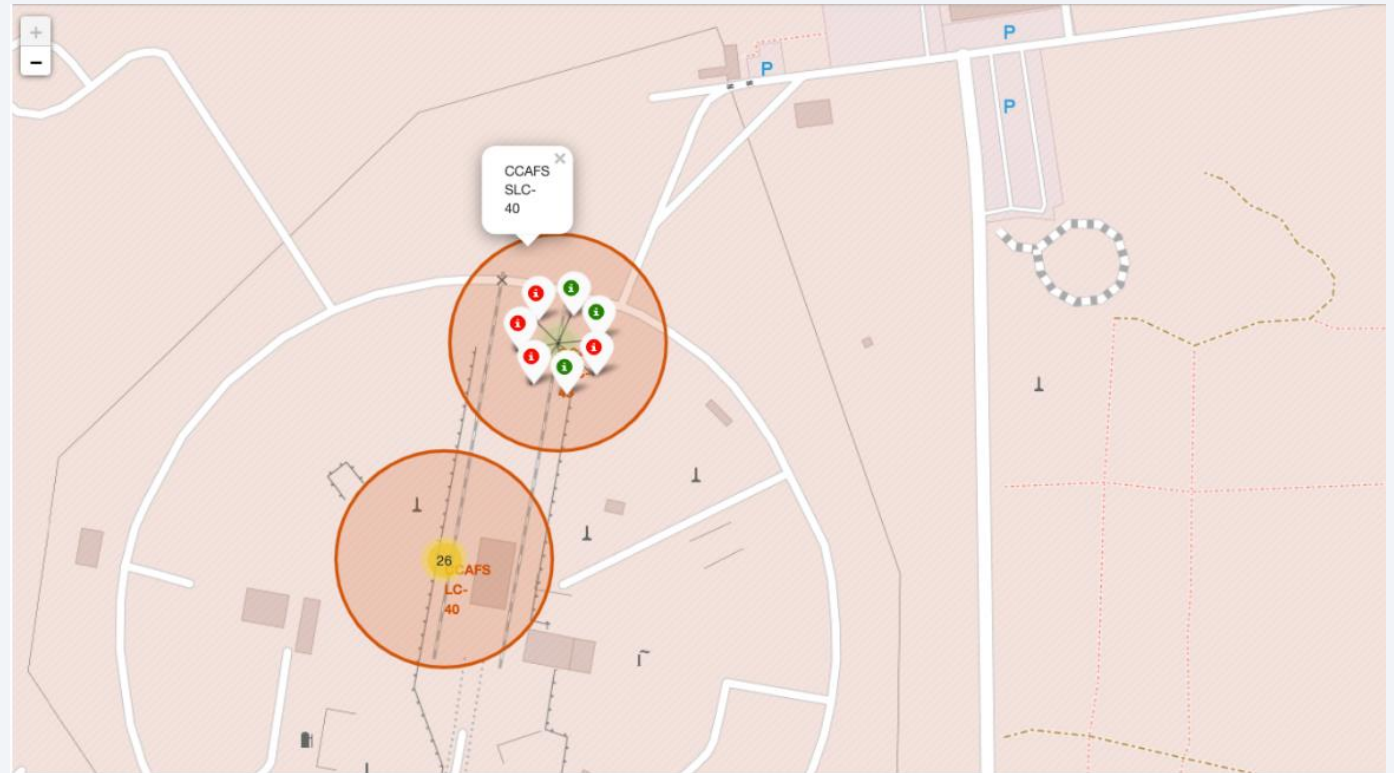# Launch Sites
# Proximities Analysis

# Map of Launch Sites

- I found that all SpaceX launch sites are located on either the West or East Cost of the United States of America
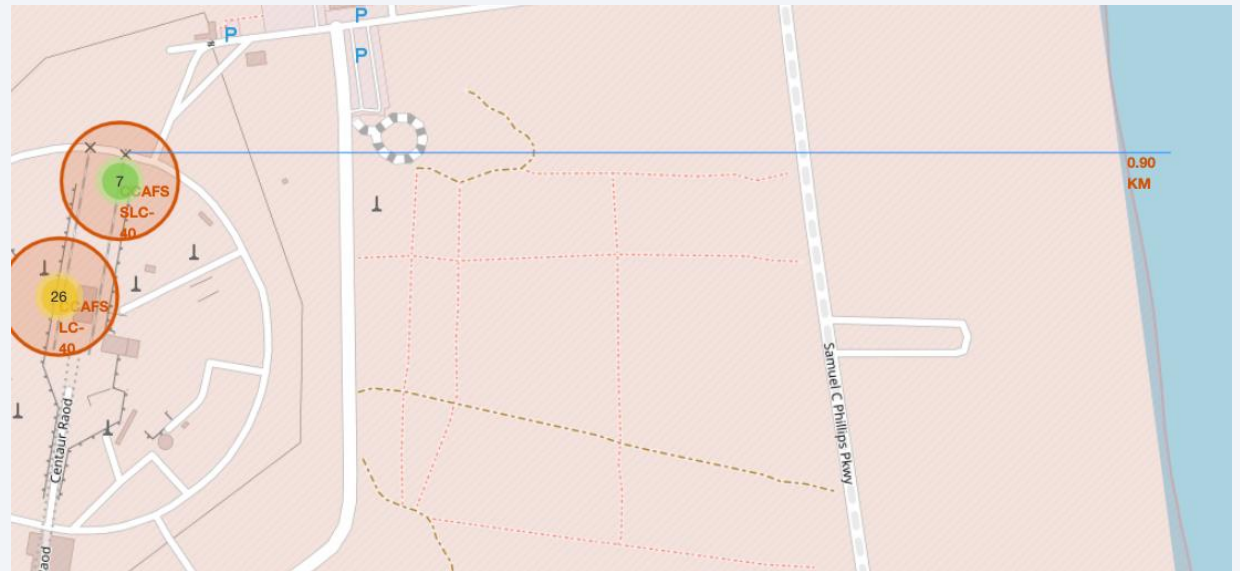
# Colored Labels Based On Landing Success

- Clustering multiple records and coloring them based on landing success can give more details on a map then if we placed all icons on the exact position of the launch

# Launch Sites and Proximities

- I found it interesting that launch sites seem to be in extremely close proximity to railroads. Most likely to transport rockets and supplies.

- I also noticed that launch sites seem to avoid crossing major highways on their short path over the ocean, in some cases less then 1 km away.
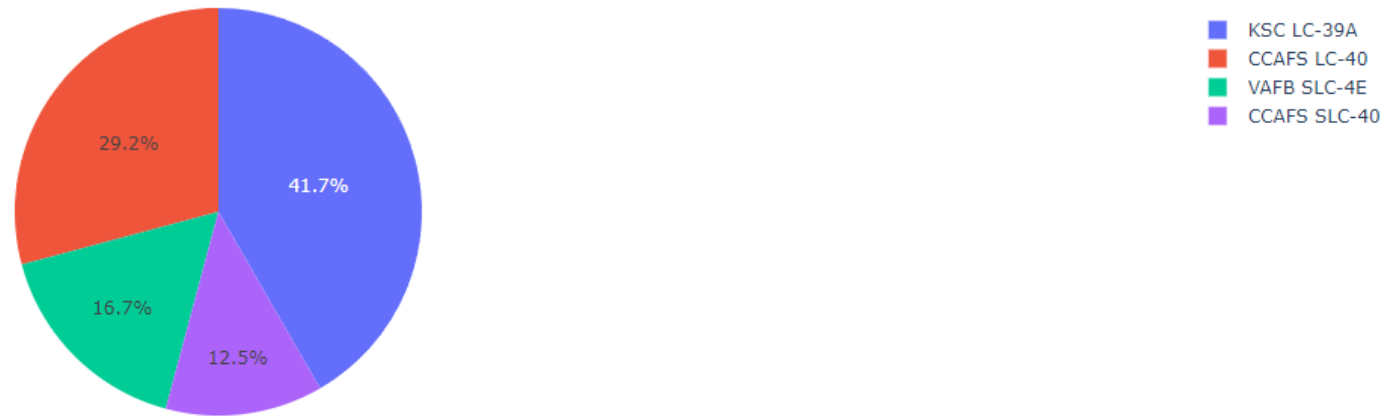
# Build a Dashboard
# with Plotly Dash

# Successful Landings for all Sites



Total Success Launches By Site

- KSC LC-39A was the most successful of all the sites

- CCAFS LC-40 was second

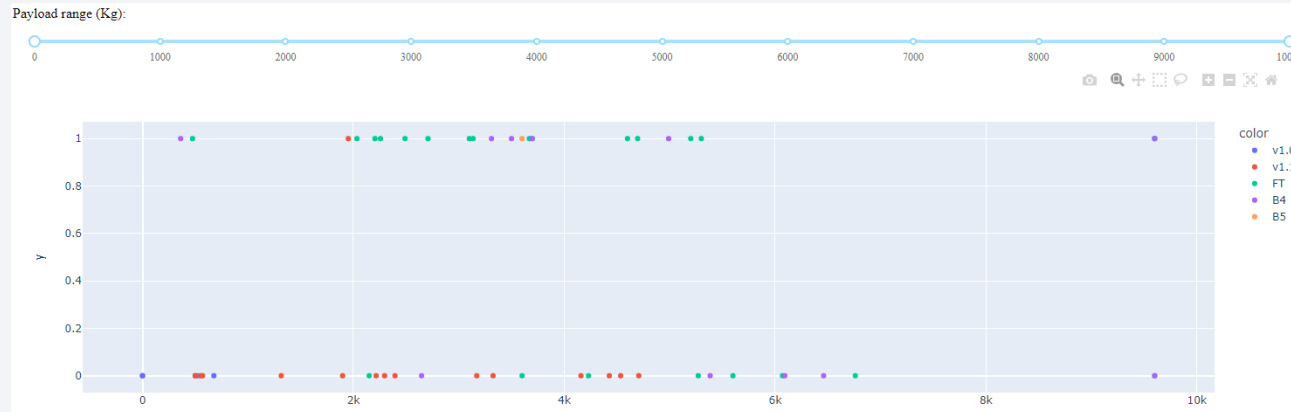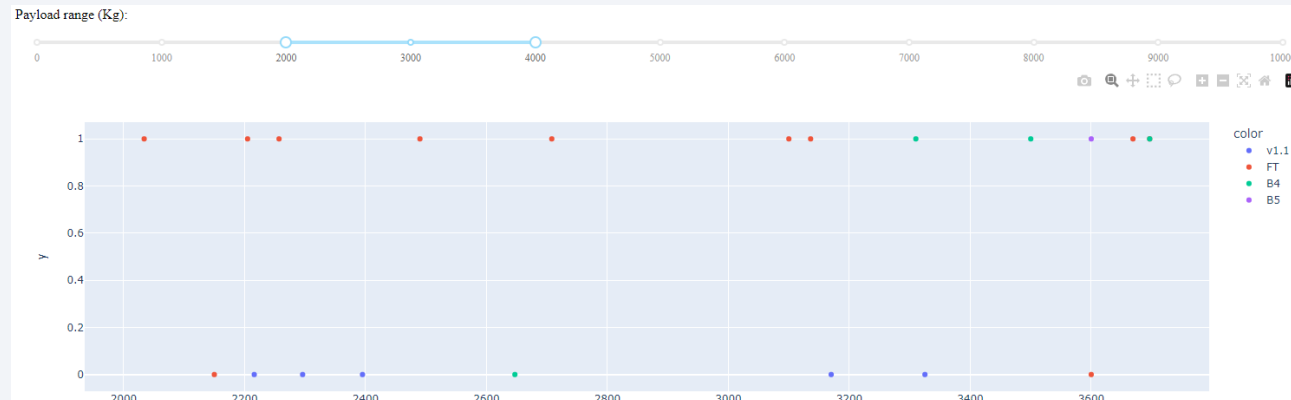- CCAFS SLC-40 was the least successful

# Most Successful Launch Site



- KSC LC-39A was the most successful launch site

- About 77% of launches ended with a successful landing

- Only about 23% of launched boosters failed to land

# Launch Success vs Payload



- Launch success for all sites and payloads
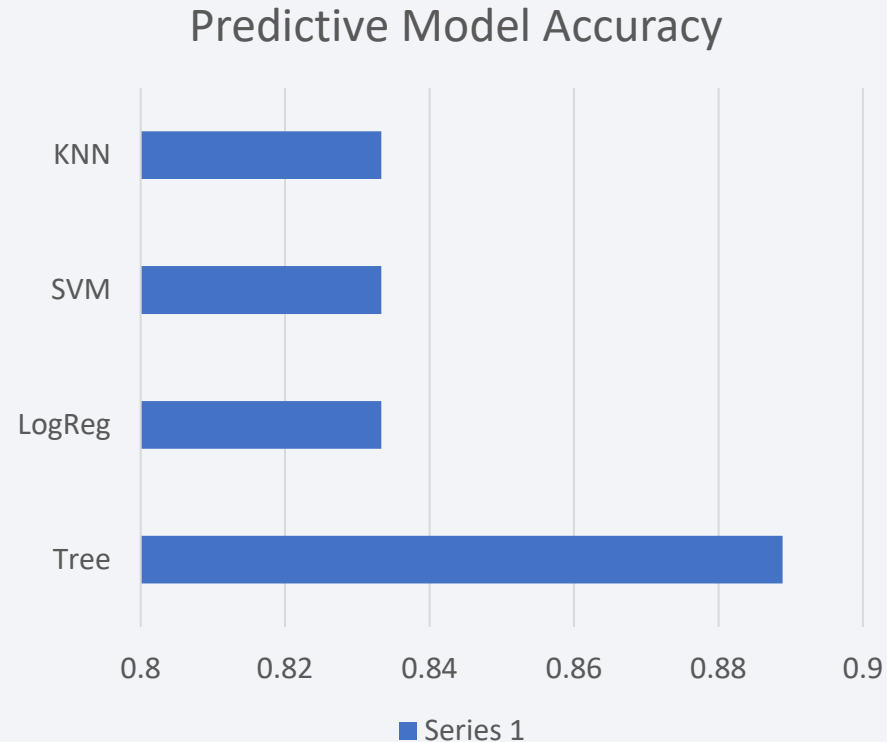
- A mix of success and failure can be seen

- Launch success for all sites and payload between 2,000 and 4,000 kg

- This was the most successful range of payloads
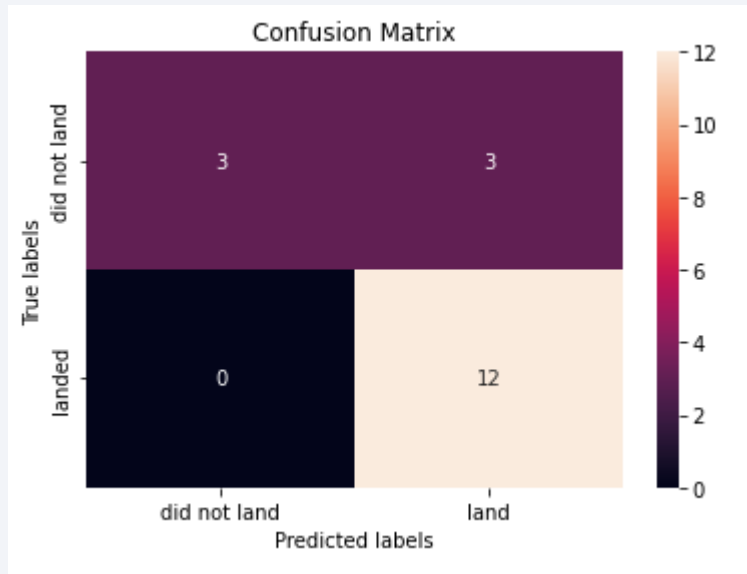
Section 5

# Predictive Analysis (Classification)

# Classification Accuracy

## Predictive Model Accuracy



- The Tree Model preformed the best when the test data was applied.

- Greater than 88% accuracy could be achieved using the tree model

# Tree Confusion Matrix



- While the tree model performed the best in my case it isn't perfect.

- The tree was good at predicting successful landings getting 12 correct and 3 incorrect.

- The tree was bad at predicting failed landings, it miss all three failures.

# Conclusions

- Successful landings can be predicted with a limited number of false positives and false negatives

- Launch sites, payload, flight number, and orbit are good predicting factors

- Predictions could improve with a larger sample of known outcomes

- Models vary in predictive accuracy

- Predictions could improve with data sets specific to the one of the predicting factors. Such as models made per Launch Site, Payload, Orbit etc.

Thank you!