**IBM Developer SKILLS NETWORK**

# Winning Space Race with Data Science

Akhil Sharma
18-Jan-2022

# Outline

- Executive Summary

- Introduction

- Methodology

- Results

- Conclusion

- Appendix

# Executive Summary

- Data Collection using API and Web Scraping

- Exploratory Data Analysis using SQL and visualization

- Interactive visual analytics using Folium and Plotly Dash

- Predictive Analysis using Classification models ( Logistic regression , SVM , Decision trees and K Nearest Neighbour )

# Introduction

In this Project we will predict if the Falcon 9 first stage will land successfully.

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.

Therefore if we can determine if the first stage will land, we can determine the cost  of a launch.

This information can be used if an alternate company wants to bid against SpaceX for a rocket launch

Section 1

# Methodology

# Methodology

## Executive Summary

- Data collection methodology:
  - SpaceX REST API
  - Web Scraping related wiki pages
- Perform data wrangling
- Perform exploratory data analysis (EDA) using visualization and SQL
- Perform interactive visual analytics using Folium and Plotly Dash
- Perform predictive analysis using classification models
  - Classification models like Logistic Regression , SVM , Decision Tree , KNN etc are applied in the Spacex Dataset. It has been trained with Train data set and then tested on test data set. Confusion Matrix has been plotted to compare the various classification models.

# Data Collection

Data Collection using API :- A  get request to the SpaceX API.

Request to the SpaceX API

Clean the data

Data Collection using Web Scraping :-Web scrap Falcon 9 launch records with BeautifulSoup:

Extracted a Falcon 9 launch records HTML table from Wikipedia

Parsed the table and converted it into a Pandas data frame

# Data Collection – SpaceX API

- GET request to read data from API.

- **GIT HUB LINK**

In [10]: static_json_url='https://cf-courses-data.s3.us.cloud-ob

We should see that the request was successfull with the 200 stat

In [11]: response.status_code

Out[11]: 200

Now we decode the response content as a Json using `.json()`

In [12]: # Use json_normalize meethod to convert the json result
data = pd.json_normalize(response.json())

# Data Collection - Scraping

- Extracted a Falcon 9 launch records HTML table from Wikipedia using BeautifulSoup object.

- **GIT HUB LINK**

# Data Wrangling

- Data Wrangling has been performed to perform Exploratory Data Analysis and determining the training Labels.

  - Calculated the number of launches on each site
  - Calculated the number and occurrence of each orbit
  - Calculated the number and occurence of mission outcome per orbit type
  - Created a landing outcome label from Outcome column

**GIT Hub URL**

# EDA with Data Visualization

Below Charts are plotted for Exploratory Data Analysis :-

> Scatter Plot to see the relationship between Flight Num & Launch Site ,
, Launch Site & Payload , Orbit & Flight Num , Payload & Orbit Type.

> Bar Chart to observe the success rate by orbit type.

> Linehart to observe the success rate by Year.

**GIT Hub URL**

# EDA with SQL

- *Display 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)*
- *List the date when the first successful landing outcome in ground pad was acheived*
- *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 failed landing_outcomes in drone ship, their booster versions, and launch site names for 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*

- *https://github.com/akhilsharma43/Firstrepo/blob/master/EDA-SQL.ipynb*

# Build an Interactive Map with Folium

- Summarize what map objects such as markers, circles, lines, etc. you created and added to a folium map
- For each launch site, added a Circle object based on its coordinate (Lat, Long) values
- Markers are added for success/failure of each launch site.
- Calculated the distance of Launch sites with its proximities.


- **https://github.com/akhilsharma43/Firstrepo/blob/master/interactive_visual_analytics_lab.ipynb**

# Build a Dashboard with Plotly Dash

- Pie chart has been added to show the success rate of various launch sites

- Slider has been added to select the payload range.

# Predictive Analysis (Classification)

- Created a column for class.
- Standardized the space X data
- Data has been splitted into Test and train data.
- Best Hyperparameter for Logistic regression , SVM , KNN and Decision tree have been selected using Grid Search
- Classifier methods has been compared based on their result on test data.

- **https://github.com/akhilsharma43/Firstrepo/blob/master/ML_prediction.ipynb**

# Results

- Exploratory data analysis results

- Interactive analytics demo in screenshots

- Predictive analysis results

Section 2

# Insights drawn from EDA

# Flight Number vs. Launch Site

scatter plot of Flight Number vs. Launch Site

```
# Plot a scatter point chart with x axis to be Flight Number and y axis to be the launch site, and hue to be the class value
sns.catplot(y="LaunchSite", x="FlightNumber", hue="Class", data=df, aspect = 5)
plt.xlabel("FlightNumber",fontsize=20)
plt.ylabel("LaunchSite",fontsize=20)
plt.show()
```



successful landing is more with less flight number except for CCAFS LS -40.

# Payload vs. Launch Site

Scatter plot of Payload vs. Launch Site



```
In [5]: # Plot a scatter point chart with x axis to be Pay Load Mass (kg) and y axis to be the launch site, and hue to be the class value
        sns.catplot(y="LaunchSite", x="PayloadMass", hue="Class", data=df, aspect = 5)
        plt.xlabel("LaunchSite",fontsize=20)
        plt.ylabel("Pay load Mass (kg)",fontsize=20)
        plt.show()
        # Plot a scatter point chart with x axis to be Flight Number and y axis to be the launch site, and hue to be the class value
```

For the VAFB-SLC launchsite there are no rockets launched for heavy payload mass(greater than 10000).

# Success Rate vs. Orbit Type

- Bar chart for the success rate of each orbit type

```
In [7]: df['SuccessRateByOrbit'] = round(100* df['Class'].groupby(df['Orbit']).transform('mean') )
        #df['Class'].groupby (by='Orbit').mean()
        # HINT use groupby method on Orbit column and get the mean of Class column
        #df['SuccessRatebyOrbit'] = df.groupby (by='Orbit').mean()['Class']
        df.head()
        plt.bar (df['Orbit'] , df['SuccessRateByOrbit'] ,color ='maroon'  )
        plt.xlabel("SuccessRateByOrbit",fontsize=20)
        plt.ylabel("Orbit",fontsize=20)
        plt.show()
```



- Orbits GEO , SSO , ES-L1 and HEo have higher success rate.

# Flight Number vs. Orbit Type

Scatter point of Flight number vs. Orbit type



```
In [8]: # Plot a scatter point chart with x axis to be FlightNumber and y axis to be the Orbit, and hue to be the class value
        sns.catplot(y="Orbit", x="FlightNumber", hue="Class", data=df, aspect = 5)
        plt.xlabel("FlightNumber",fontsize=20)
        plt.ylabel("Orbit",fontsize=20)
        plt.show()
```

LEO orbit the Success appears related to the number of flights; on the other hand, there seems to be no relationship between flight number when in GTO orbit.

# Payload vs. Orbit Type

- Scatter point of payload vs. orbit type



```
# Plot a scatter point chart with x axis to be Payload and y axis to be the Orbit, and hue to be the class value
sns.catplot(y="Orbit", x="PayloadMass", hue="Class", data=df, aspect = 5)
plt.xlabel("Payload",fontsize=20)
plt.ylabel("Orbit",fontsize=20)
plt.show()
```

With heavy payloads the successful landing or positive landing rate are more for Polar,LEO and ISS.However for GTO we cannot distinguish this well as both positive landing rate and negative landing(unsuccessful mission) are both there here.

# Launch Success Yearly Trend

- Line chart of yearly average success rate



```
In [36]: # Plot a line chart with x axis to be the extracted year and y axis to be the success rate
         plt.plot(year, df['successbyyear'])
         plt.show()
```

Success rate since 2013 kept increasing till 2020

# All Launch Site Names

- Find the names of the unique launch sites



select distinct launch_site from TQW1863...    Run time: **0.007 s**

**Result set 1**    Q Find

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

# Launch Site Names Begin with 'CCA'

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

# Total Payload Mass

- Calculate the total payload carried by boosters from NASA

# Average Payload Mass by F9 v1.1

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

# First Successful Ground Landing Date

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

# 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

# Total Number of Successful and Failure Mission Outcomes

- Calculate the total number of successful and failure mission outcomes

# Boosters Carried Maximum Payload

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

# 2015 Launch Records

- List the failed landing_outcomes in drone ship, their booster versions, and launch site names for in year 2015

# 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 date , landing__outcome , rank ( ) over
( partition by landing__outcome order by date )
from TQW18636.SPACEX where date between '2010-06-04' and '2017-03-20'
```

| DATE | LANDING__OUTCOME | 3 |
|------|------------------|---|
| 2015-11-02 | Controlled (ocean) | 1 |
| 2015-10-01 | Failure (drone ship) | 1 |
| 2016-04-03 | Failure (drone ship) | 2 |
| 2010-08-12 | Failure (parachute) | 1 |
| 2012-08-10 | No attempt | 1 |

Result set is truncated, only the first 15 rows have been loaded. Select "View all loaded data" on the right top of the result to view all loaded rows.
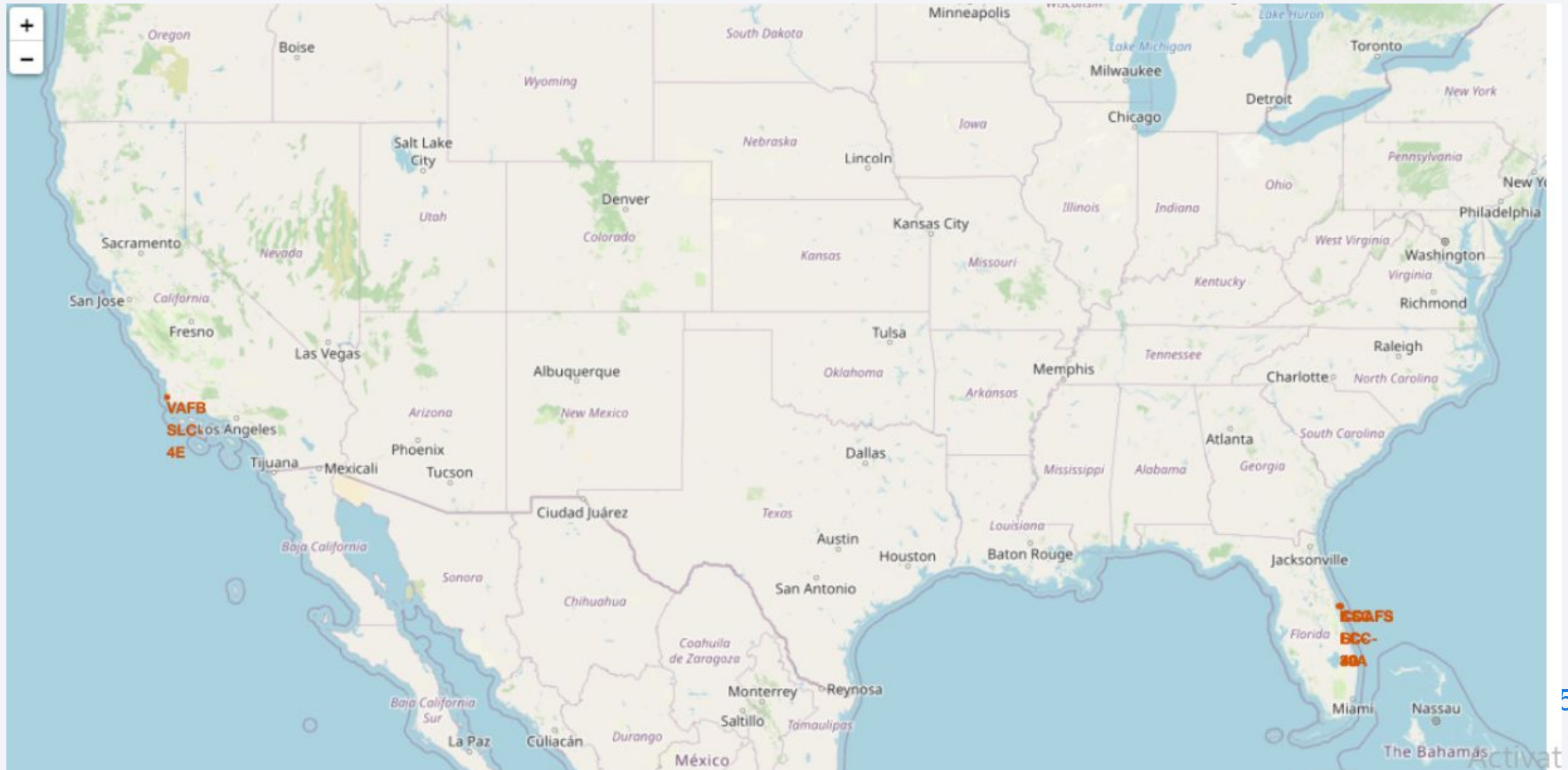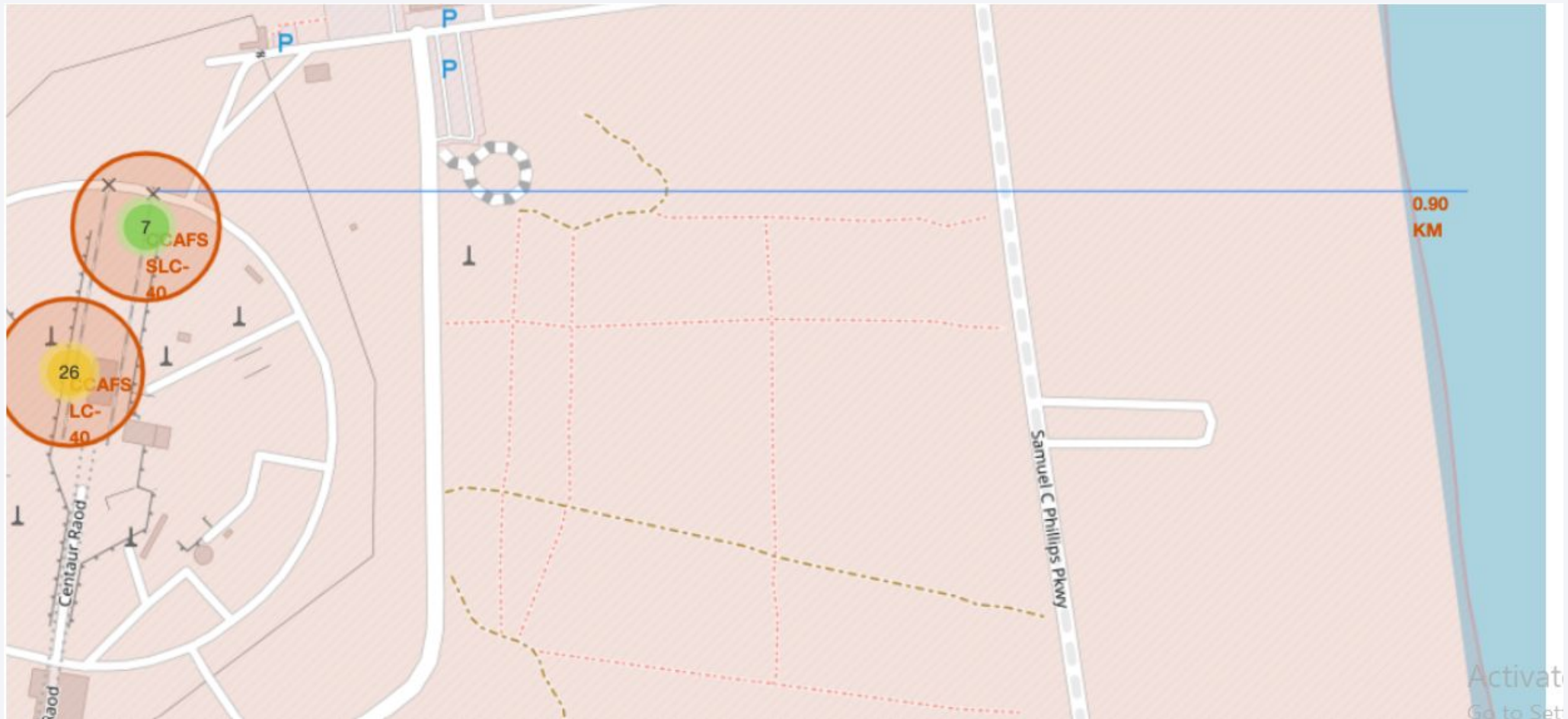
More

Section 4

# Launch Sites Proximities Analysis

# Launch Site Locations

# Success/Failure at Launch Sites

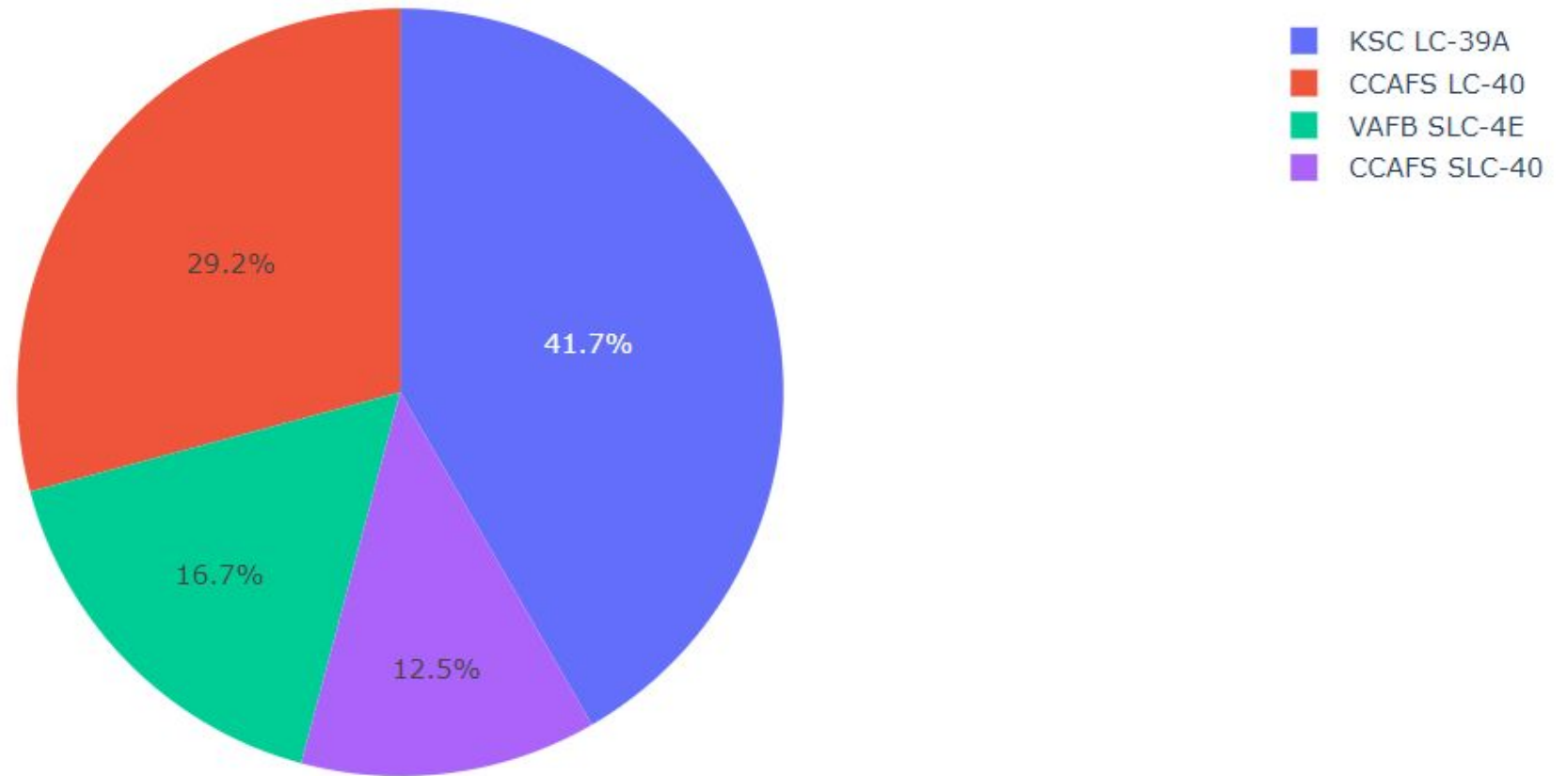# Distance of Proximities from Launch Site

Section 5

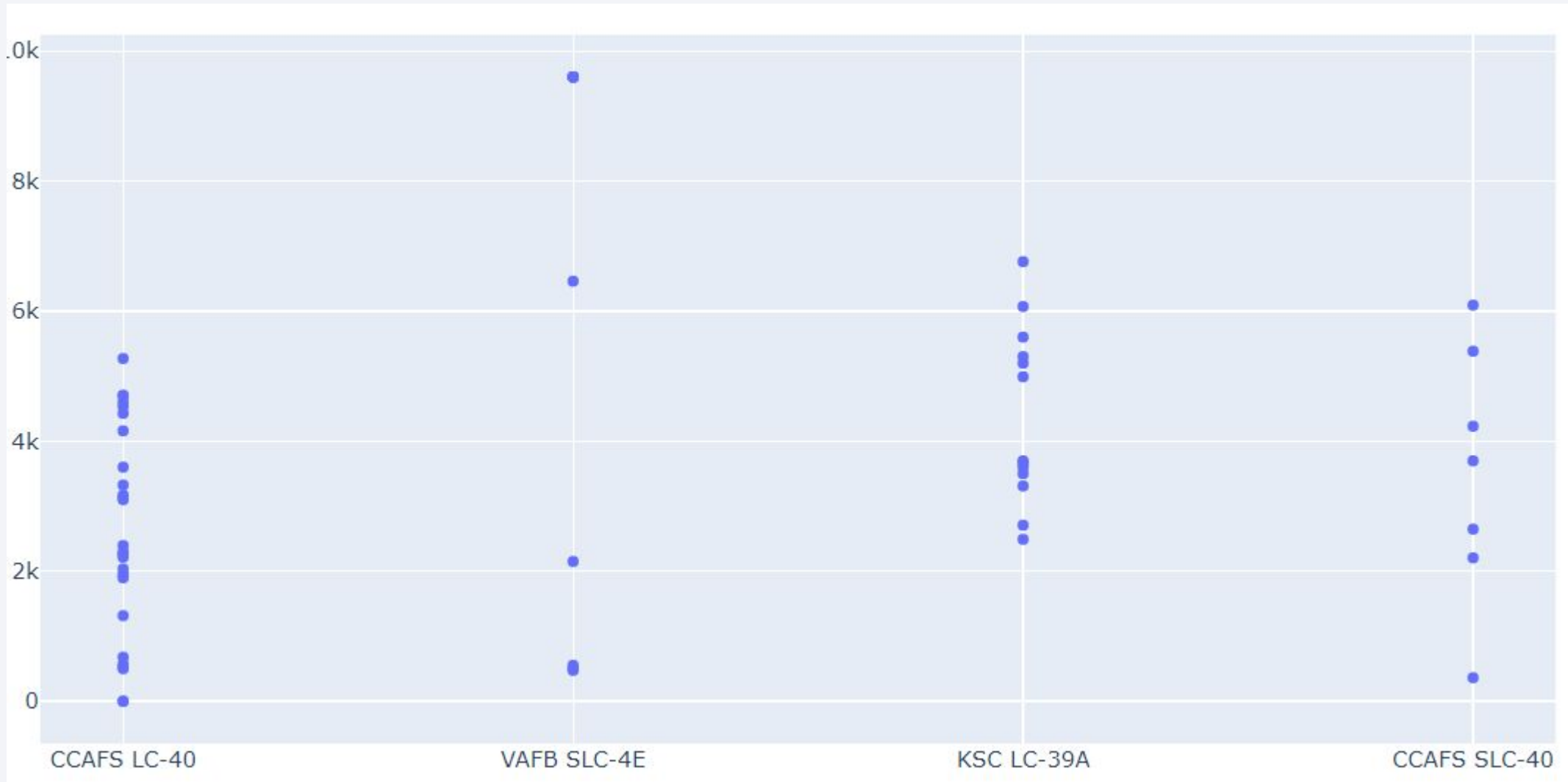# Build a Dashboard with Plotly Dash

# Success Pie Chart



success-pie-chart

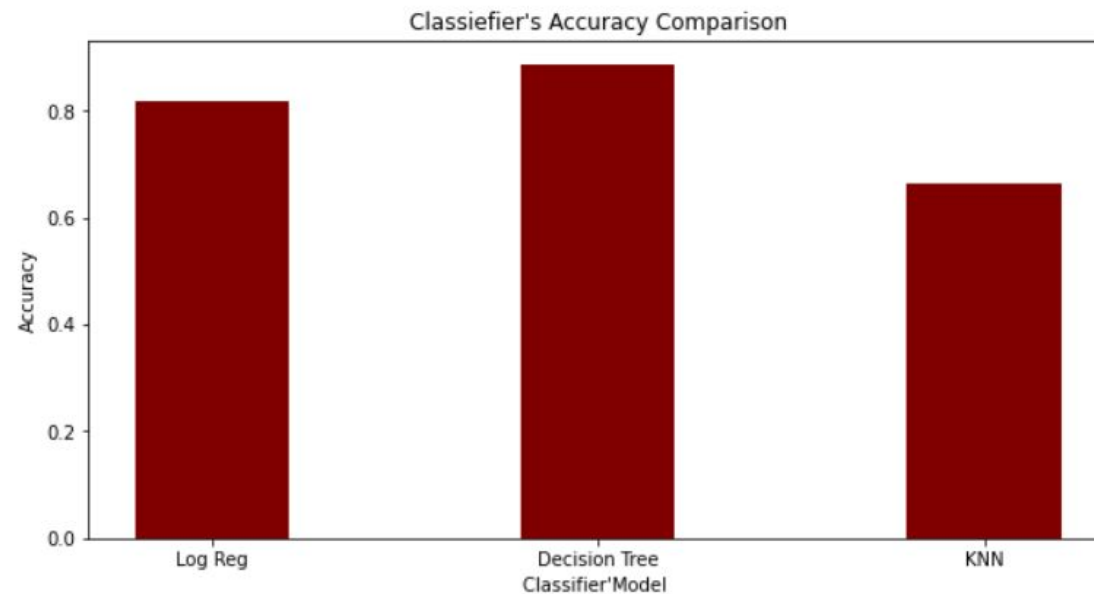# <Dashboard Screenshot 2>

Section 6

# Predictive Analysis (Classification)

# Classification Accuracy

Visualize the built model accuracy for all built classification models, in a bar chart

```
# creating the bar plot
plt.bar(['Log Reg','Decision Tree','KNN'], [logreg_cv.best_score_,tree_cv.best_score_,knn_cv.best_score_], color ='maroon',
        width = 0.4)

plt.xlabel("Classifier'Model ")
plt.ylabel("Accuracy")
plt.title("Classiefier's Accuracy Comparison")
plt.show()
```
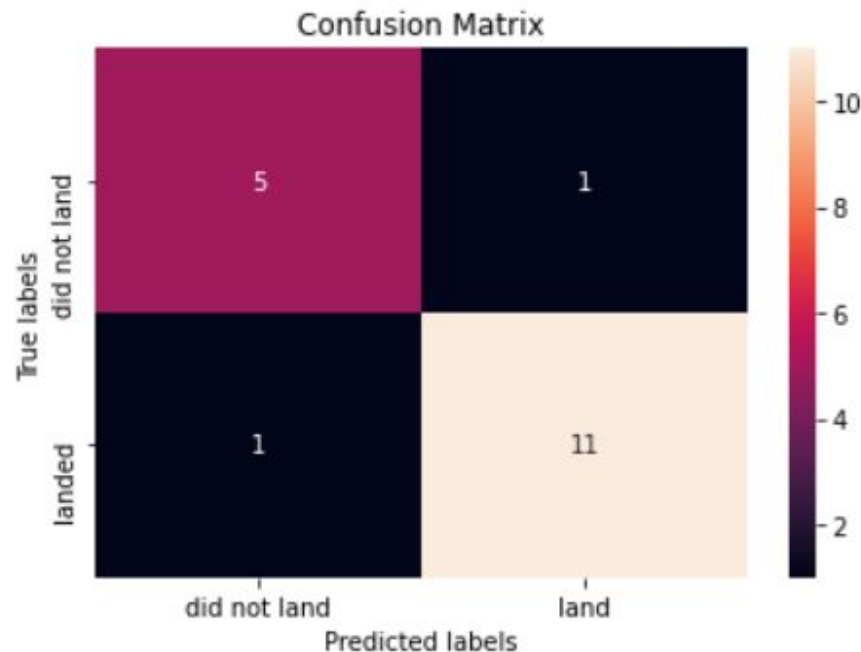
Classiefier's Accuracy Comparison

Decision Tree has the highest Accuracy among all classifier's models/

# Confusion Matrix

- Confusion Matrix for Decision Tree. 1 outcome predicted by model as Landed was actually not landed. Similarly 1 outcome predicted as not landed is actually landed.

# Conclusions

- Decision Tree performed better than other classifier models in Spacex Data set.

- Its Accuracy is around 89%.

- 1 outcome is False positive and 1 outcome is False Negative out of total 18 outcomes.

Thank you!