# Space X Falcon

# The First Stage

**Deepak Kumar Singh**
**07/23/2023**

# Contents

2023

# Executive Summary

## Summary of Methodologies

The research attempts to identify the factors for a successful rocket landing. To make this determination, the following methodologies where used:

- **Collect** data using SpaceX REST API and web scraping techniques
- **Wrangle** data to create success/fail outcome variable
- **Explore** data with data visualization techniques, considering the following factors: payload, launch site, flight number and yearly trend
- **Analyze** the data with SQL, calculating the following statistics: total payload, payload range for successful launches, and total #of successful and failed outcomes
- **Explore** launch site success rates and proximity to geographical markers
- **Visualize** the launch sites with the most success and successful payload ranges
- **Build Models** to predict landing outcomes using logistic regression, support vector machine (SVM), decision tree and K-nearest neighbor (KNN)

## Results

### Exploratory Data Analysis:
- Launch success has improved over time
- KSC LC-39A has the highest success rate among landing sites
- Orbits ES-L1, GEO, HEO, and SSO have a 100% success rate

### Predictive Analytics:
- All models performed similarly on the test set. The decision tree model slightly outperformed

### Visualization/Analytics:
- Most launch sites are near the equator, and all are close to the coast

# Introduction

## Background

Space X 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 Space X 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 space X for a rocket launch. This goal of the project is to create a machine learning pipeline to predict if the first stage will land successfully.

## Problems you want to find answers

• How payload mass, launch site, number of flights, and orbits affect first-stage landing success
• Rate of successful landings over time
• Best predictive model for successful landing (binary classification)

Section -1
Methodology

# Methodology

## Steps

- **Collect** data using SpaceX REST API and web scraping techniques

- **Wrangle** data – by filtering the data, handling missing values and applying one hot encoding – to prepare the data for analysis and modeling

- **Explore** data via EDA with SQL and data visualization techniques

- **Visualize** the data using Folium and Plotly Dash

- **Build Models** to predict landing outcomes using classification models. Tune and evaluate models to find best model and parameters

# Data Collection –

❑ The data was collected using various methods

   ❑ Data collection was done using get request to the SpaceX API.

   ❑ Next, we decoded the response content as a Json using .json() function call and turn it into a pandas dataframe using .json_normalize().

   ❑ We then cleaned the data, checked for missing values and fill in missing values where necessary.

   ❑ In addition, we performed web scraping from Wikipedia for Falcon 9 launch records with BeautifulSoup.

   ❑ The objective was to extract the launch records as HTML table, parse the table and convert it to a pandas dataframe for future analysis.

# Data Collection – API

We used the get request to the SpaceX API to collect data, clean the requested data and did some basic data wrangling and formatting.

The link to the notebook is:

jupyter-labs-spacex-data-collection-api.ipynb

1.  Get request for rocket launch data using API

```
In [6]:  spacex_url="https://api.spacexdata.com/v4/launches/past"
```

```
In [7]:  response = requests.get(spacex_url)
```

2.  Use json_normalize method to convert json result to dataframe

```
In [12]:  # Use json_normalize method to convert the json result into a dataframe

          # decode response content as json
          static_json_df = res.json()
```

```
In [13]:  # apply json_normalize
          data = pd.json_normalize(static_json_df)
```

3.  We then performed data cleaning and filling in the missing values

```
In [30]:  rows = data_falcon9['PayloadMass'].values.tolist()[0]

          df_rows = pd.DataFrame(rows)
          df_rows = df_rows.replace(np.nan, PayloadMass)

          data_falcon9['PayloadMass'][0] = df_rows.values
          data_falcon9
```

# Data Collection – Web Scraping



We applied web scrapping to webscrap Falcon 9 launch records with BeautifulSoup

We parsed the table and converted it into a pandas dataframe.

The link to the notebook is:

https://github.com/Deepaksingh2310/SpaceX-Capstone-project/blob/main/jupyter-labs-webscraping%20(1).ipynb

# Data Wrangling

o  We performed exploratory data analysis and determined the training labels.

o  We calculated the number of launches at each site, and the number and occurrence of each orbits

o  We created landing outcome label from outcome column and exported the results to csv.

o  The link to the notebook is

o  https://github.com/Deepaksingh2310/SpaceX-Capstone-project/blob/main/labs-jupyter-spacex-data_wrangling_jupyterlite.jupyterlite.ipynb

# EDA with Visualization

We explored the data by visualizing the relationship between flight number and launch Site, payload and launch site, success rate of each orbit type, flight number and orbit type, the launch success yearly trend.

The link to the notebook is
https://github.com/Deepaksingh2310/SpaceX-Capstone-project/blob/main/jupyter-labs-eda-dataviz.ipynb.jupyterlite.ipynb

2023



Plot of success rate by class of each Orbits



Plot of launch success yearly trend

# EDA with SQL

• We loaded the SpaceX dataset into a PostgreSQL database without leaving the jupyter notebook.

• We applied EDA with SQL to get insight from the data. We wrote queries to find out for instance:

   • The names of unique launch sites in the space mission.

   • The total payload mass carried by boosters launched by NASA (CRS)

   • The average payload mass carried by booster version F9 v1.1

   • The total number of successful and failure mission outcomes

   • The failed landing outcomes in drone ship, their booster version and launch site names.

• The link to the notebook is  https://github.com/Deepaksingh2310/SpaceX-Capstone-project/blob/main/jupyter-labs-eda-sql-coursera_sqllite.ipynb

# Map with Folium

•We marked all launch sites, and added map objects such as markers, circles, lines to mark the success or failure of launches for each site on the folium map.

•We assigned the feature launch outcomes (failure or success) to class 0 and 1.i.e., 0 for failure, and 1 for success.

•Using the color-labeled marker clusters, we identified which launch sites have relatively high success rate.

•We calculated the distances between a launch site to its proximities. We answered some question for instance:

    -Are launch sites near railways, highways and coastlines.

    -Do launch sites keep certain distance away from cities.

# Dashboard with Plotly Dash

- We built an interactive dashboard with Plotly dash

- We plotted pie charts showing the total launches by a certain sites

- We plotted scatter graph showing the relationship with Outcome and Payload Mass (Kg) for the different booster version.

- The link to the notebook is

- https://github.com/Deepaksingh2310/SpaceX-Capstone-project/blob/main/spacex_dash_app.py

# Predictive Analytics

- We loaded the data using numpy and pandas, transformed the data, split our data into training and testing.

- We built different machine learning models and tune different hyperparameters using GridSearchCV.

- We used accuracy as the metric for our model, improved the model using feature engineering and algorithm tuning.

- We found the best performing classification model.

- The link to the notebook is

- https://github.com/Deepaksingh2310/SpaceX-Capstone-project/blob/main/SpaceX_Machine_Learning_Prediction_Part_5.jupyterlite.ipynb

# Section 2
Results

# Results Summary

**Exploratory Data Analysis**

• Launch success has improved over time

• KSC LC-39A has the highest success rate among landing sites

• Orbits ES-L1, GEO, HEO and SSO have a 100% success rate

**Visual Analytics**

• Most launch sites are near the equator, and all are close to the coast

• Launch sites are far enough away from anything a failed launch can damage (city, highway, railway), while still close enough to bring people and material to support launch activities

**Predictive Analytics**

• Decision Tree model is the best predictive model for the dataset

# Flight Number vs Launch Site

**Exploratory Data Analysis**

• **Earlier flights** had a **lower success rate** (**blue = fail**)

• **Later flights** had a **higher success rate** (**orange = success**)

• Around half of launches were from CCAFS SLC 40 launch site

• VAFB SLC 4E and KSC LC 39A have higher success rates

• We can infer that new launches have a higher success rate

# Payload vs. Launch Site

**Exploratory Data Analysis**

- Typically, the **higher** the **payload mass** (kg), the **higher** the **success rate**

- Most launces with a payload greater than 7,000 kg were successful

- KSC LC 39A has a 100% success rate for launches less than 5,500 kg

- VAFB SKC 4E has not launched anything greater than ~10,000 kg

# Success Rate by Orbit

**Exploratory Data Analysis**

- **100% Success Rate**: ES-L1, GEO, HEO and SSO

- **50%-80% Success Rate**: GTO, ISS, LEO, MEO, PO

- **0% Success Rate**: SO



Plot of success rate by class of each Orbits

# Flight Number vs. Orbit

**Exploratory Data Analysis**

- The success rate typically increases with the number of flights for each orbit

- This relationship is highly apparent for the LEO orbit

- The GTO orbit, however, does not follow this trend

# Payload vs. Orbit

**Exploratory Data Analysis**

• Heavy payloads are better with LEO, ISS and PO orbits

• The GTO orbit has mixed success with heavier payloads

# Launch Success over Time

**Exploratory Data Analysis**

- The success rate improved from 2013-2017 and 2018-2019

- The success rate decreased from 2017-2018 and from 2019-2020

- Overall, the success rate has improved since 2013

# Launch Site Information

**Launch Site Names**

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

**Records with Launch Site Starting with CCA**

- Displaying 5 records below

Display the names of the unique launch sites in the space mission

```
In [10]:   task_1 = '''
              SELECT DISTINCT LaunchSite
              FROM SpaceX
           '''
           create_pandas_df(task_1, database=conn)
```

Out[10]:

| | launchsite |
|---|---|
| 0 | KSC LC-39A |
| 1 | CCAFS LC-40 |
| 2 | CCAFS SLC-40 |
| 3 | VAFB SLC-4E |

Display 5 records where launch sites begin with the string 'CCA'

```
In [11]:   task_2 = '''
              SELECT *
              FROM SpaceX
              WHERE LaunchSite LIKE 'CCA%'
              LIMIT 5
           '''
           create_pandas_df(task_2, database=conn)
```

Out[11]:

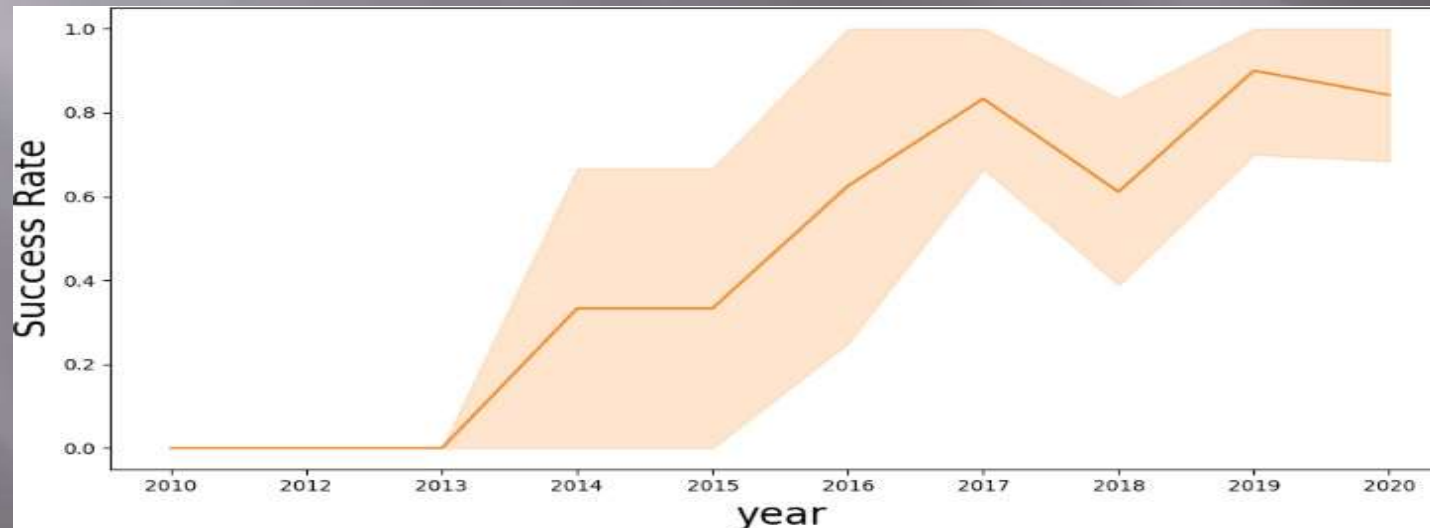| | date | time | boosterversion | launchsite | payload | payloadmasskg | orbit | customer | missionoutcome | landingoutcome |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 2010-04-06 | 18:45:00 | F9 v1.0 B0003 | CCAFS LC-40 | Dragon Spacecraft Qualification Unit | 0 | LEO | SpaceX | Success | Failure (parachute) |
| 1 | 2010-08-12 | 15:43:00 | F9 v1.0 B0004 | CCAFS LC-40 | Dragon demo flight C1, two CubeSats, barrel of... | 0 | LEO (ISS) | NASA (COTS) NRO | Success | Failure (parachute) |
| 2 | 2012-05-22 | 07:44:00 | F9 v1.0 B0005 | CCAFS LC-40 | Dragon demo flight C2 | 525 | LEO (ISS) | NASA (COTS) | Success | No attempt |
| 3 | 2012-08-10 | 00:35:00 | F9 v1.0 B0006 | CCAFS LC-40 | SpaceX CRS-1 | 500 | LEO (ISS) | NASA (CRS) | Success | No attempt |
| 4 | 2013-01-03 | 15:10:00 | F9 v1.0 B0007 | CCAFS LC-40 | SpaceX CRS-2 | 677 | LEO (ISS) | NASA (CRS) | Success | No attempt |

# Payload Mass

**Total Payload Mass**

- **45,596 kg** (total) carried by boosters launched by NASA (CRS)

```
%sql SELECT SUM(PAYLOAD_MASS__KG_) \
    FROM SPACEXTBL \
    WHERE CUSTOMER = 'NASA (CRS)';

 * ibm_db_sa://yyy33800:***@1bbf73c5-d84a-4
   sqlite:///my_data1.db
Done.

    1

 45596
```

**Average Payload Mass**

- **2,928 kg** (average) carried by booster version F9 v1.1

```
%sql SELECT AVG(PAYLOAD_MASS__KG_) \
    FROM SPACEXTBL \
    WHERE BOOSTER_VERSION = 'F9 v1.1';

 * ibm_db_sa://yyy33800:***@1bbf73c5-d84a-4
   sqlite:///my_data1.db
Done.

    1

 2928
```

# Landing & Mission Info

We observed that the dates of the first successful landing outcome on ground pad was 22nd December 2015

```
In [14]:  task_5 = '''
              SELECT MIN(Date) AS FirstSuccessfull_landing_date
              FROM SpaceX
              WHERE LandingOutcome LIKE 'Success (ground pad)'
              '''
          create_pandas_df(task_5, database=conn)
```

```
Out[14]:     firstsuccessfull_landing_date

         0              2015-12-22
```

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

We used the WHERE clause to filter for boosters which have successfully landed on drone ship and applied the AND condition to determine successful landing with payload mass greater than 4000 but less than 6000

```
In [15]:  task_6 = '''
              SELECT BoosterVersion
              FROM SpaceX
              WHERE LandingOutcome = 'Success (drone ship)'
                  AND PayloadMassKG > 4000
                  AND PayloadMassKG < 6000
              '''
          create_pandas_df(task_6, database=conn)
```

```
Out[15]:     boosterversion

         0     F9 FT B1022

         1     F9 FT B1026

         2     F9 FT B1021.2

         3     F9 FT B1031.2
```

# Boosters

**Carrying Max Payload**

- 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

List the names of the booster_versions which have carried the maximum payload mass. Use a subquery

```
In [17]:  task_8 = '''
            SELECT BoosterVersion, PayloadMassKG
            FROM SpaceX
            WHERE PayloadMassKG = (
                                    SELECT MAX(PayloadMassKG)
                                    FROM SpaceX
                                    )
            ORDER BY BoosterVersion
            '''
          create_pandas_df(task_8, database=conn)
```

Out[17]:

|     | boosterversion | payloadmasskg |
| --- | -------------- | ------------- |
| 0   | F9 B5 B1048.4  | 15600         |
| 1   | F9 B5 B1048.5  | 15600         |
| 2   | F9 B5 B1049.4  | 15600         |
| 3   | F9 B5 B1049.5  | 15600         |
| 4   | F9 B5 B1049.7  | 15600         |
| 5   | F9 B5 B1051.3  | 15600         |
| 6   | F9 B5 B1051.4  | 15600         |
| 7   | F9 B5 B1051.6  | 15600         |
| 8   | F9 B5 B1056.4  | 15600         |
| 9   | F9 B5 B1058.3  | 15600         |
| 10  | F9 B5 B1060.2  | 15600         |
| 11  | F9 B5 B1060.3  | 15600         |

# Failed Landings on Drone Ship

**In 2015**

- Showing month, date, booster version, launch site and landing outcome

```
%sql SELECT substr(Date,4,2) as month, DATE,BOOSTER_VERSION, LAUNCH_SITE, [Landing _Outcome] \
FROM SPACEXTBL \
where [Landing _Outcome] = 'Failure (drone ship)' and substr(Date,7,4)='2015';
```

* sqlite:///my_data1.db
Done.

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

# Count of Successful Landings

**Ranked Descending**

- Count of landing outcomes between 2010-06-04 and 2017-03-20 in descending order

```
%sql SELECT [Landing _Outcome], count(*) as count_outcomes \
FROM SPACEXTBL \
WHERE DATE between '04-06-2010' and '20-03-2017' group by [Landing _Outcome] order by count_outcomes DESC;
```

* sqlite:///my_data1.db
Done.

| Landing _Outcome | count_outcomes |
|---|---|
| Success | 20 |
| No attempt | 10 |
| Success (drone ship) | 8 |
| Success (ground pad) | 6 |
| Failure (drone ship) | 4 |
| Failure | 3 |
| Controlled (ocean) | 3 |
| Failure (parachute) | 2 |
| No attempt | 1 |

Section 4
Launch Site
Analysis

# Launch Sites

**With Markers**

- **Near Equator**: the closer the launch site to the equator, the **easier** it is **to launch** to equatorial orbit, and the more help you get from Earth's rotation for a prograde orbit. Rockets launched from sites near the equator get an **additional natural boost** - due to the rotational speed of earth - that **helps save the cost** of putting in extra fuel and boosters.
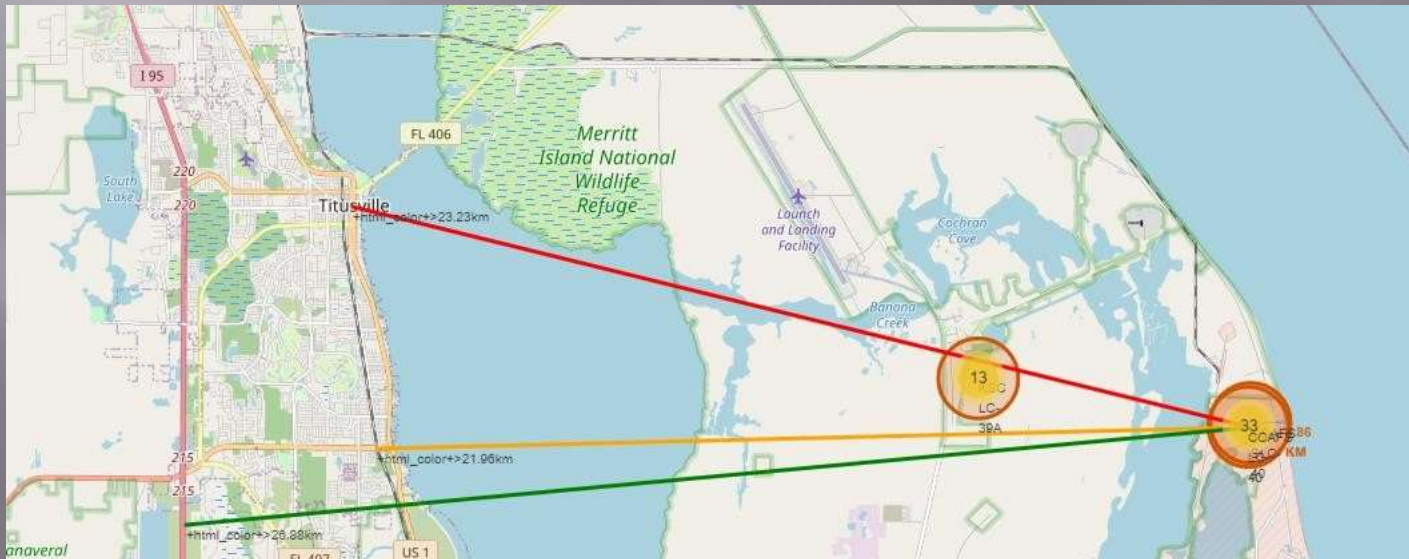
# Launch Outcomes

**At Each Launch Site**

- **Outcomes**:

- **Green** markers for successful launches

- **Red** markers for unsuccessful launches

- Launch site **CCAFS SLC-40** has a **3/7 success rate** (**42.9%**)

# Distance to Proximities

**CCAFS SLC-40**

- **.86 km** from nearest coastline

- **21.96 km** from nearest railway

- **23.23 km** from nearest city
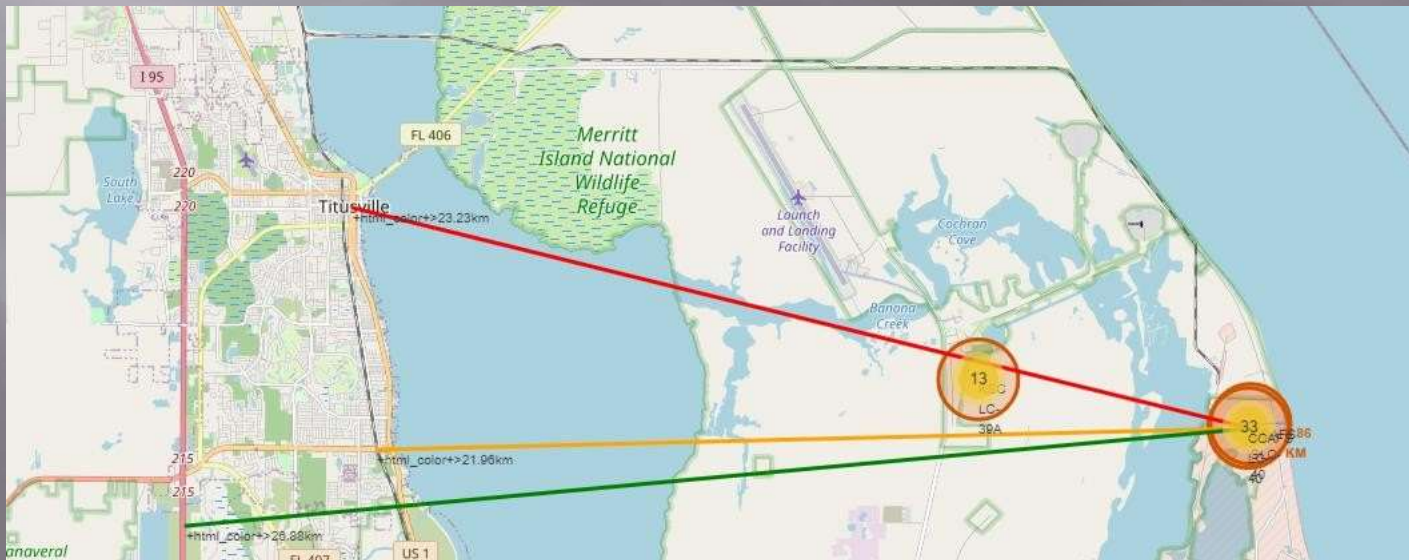
- **26.88 km** from nearest highway

# Distance to Proximities

**CCAFS SLC-40**

- **Coasts**: help ensure that spent stages dropped along the launch path or failed launches don't fall on people or property.

- **Safety / Security:** needs to be an exclusion zone around the launch site to keep unauthorized people away and keep people safe.

- **Transportation/Infrastructure and Cities**: need to be away from anything a failed launch can damage, but still close enough to roads/rails/docks to be   able to bring people and material to or from it in support of launch activities.
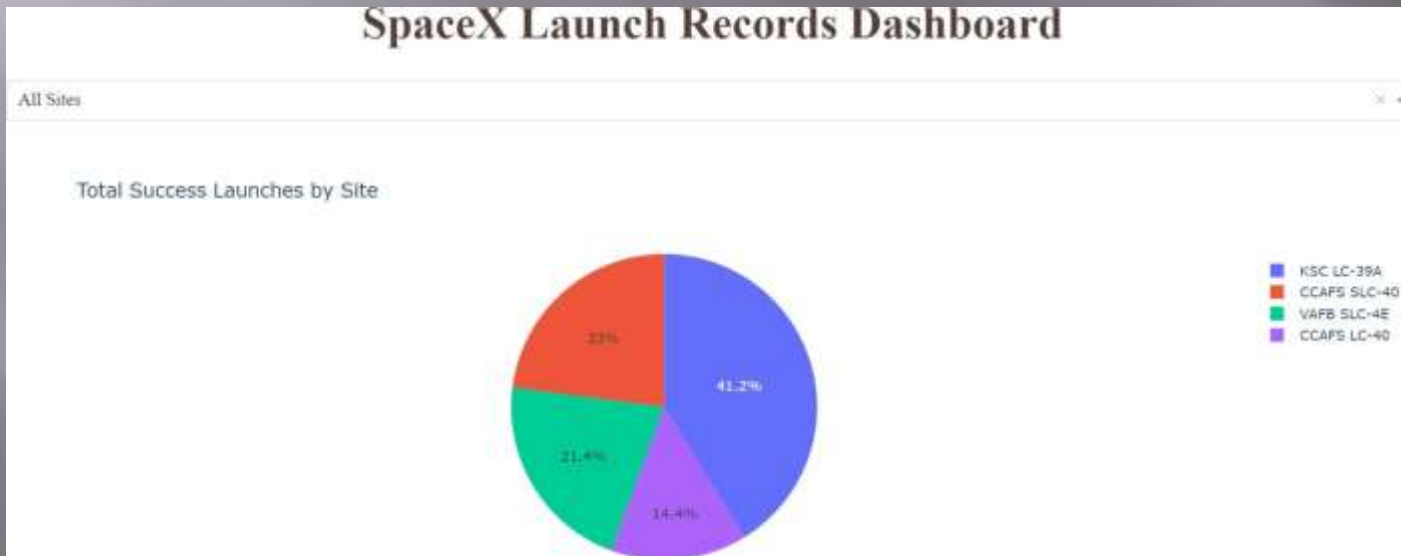
Section 5

Dashboard with  Plotly

# Launch Success by Site

**Success as Percent of Total**
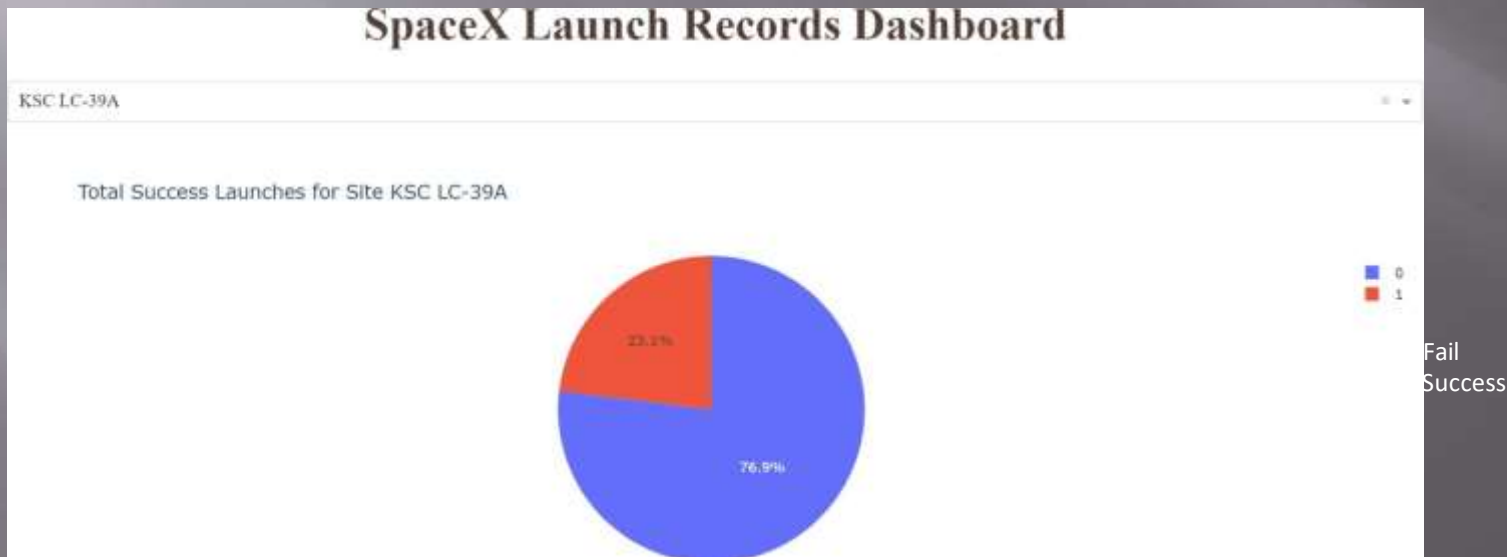
- **KSC LC-39A** has the **most successful launches** amongst launch sites (**41.2%**)



SpaceX Launch Records Dashboard

All Sites

Total Success Launches by Site

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

# Launch Success (KSC LC-29A)

**Success as Percent of Total**

- **KSC LC-39A** has the **highest success rate** amongst launch sites (**76.9%**)

- 10 successful launches and 3 failed launches

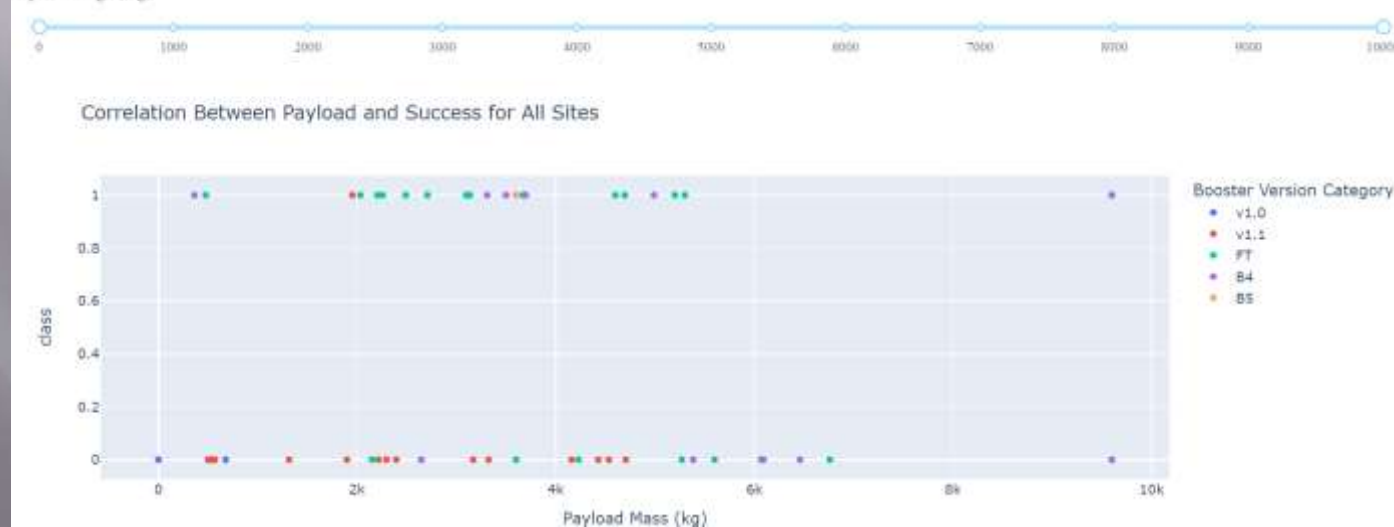# Payload Mass and Success

**By Booster Version**

- **Payloads between 2,000 kg** and **5,000 kg** have the **highest success rate**

- 1 indicating successful outcome and 0 indicating an unsuccessful outcome

# Section 5
# Predictive  Analytics

# Classification

**Accuracy**

- **All** the **models** performed at about the same level and had the **same scores** and **accuracy.** This is likely due to the **small dataset**. The **Decision Tree model slightly outperformed** the rest when looking at .best_score_

- .best_score_ is the average of all cv folds for a single combination of the parameters

| | LogReg | SVM | Tree | KNN |
|---|---|---|---|---|
| Jaccard_Score | 0.800000 | 0.800000 | 0.800000 | 0.800000 |
| F1_Score | 0.888889 | 0.888889 | 0.888889 | 0.888889 |
| Accuracy | 0.833333 | 0.833333 | 0.833333 | 0.833333 |

```python
models = {'KNeighbors':knn_cv.best_score_,
          'DecisionTree':tree_cv.best_score_,
          'LogisticRegression':logreg_cv.best_score_,
          'SupportVector': svm_cv.best_score_}

bestalgorithm = max(models, key=models.get)
print('Best model is', bestalgorithm,'with a score of', models[bestalgorithm])
if bestalgorithm == 'DecisionTree':
    print('Best params is :', tree_cv.best_params_)
if bestalgorithm == 'KNeighbors':
    print('Best params is :', knn_cv.best_params_)
if bestalgorithm == 'LogisticRegression':
    print('Best params is :', logreg_cv.best_params_)
if bestalgorithm == 'SupportVector':
    print('Best params is :', svm_cv.best_params_)

Best model is DecisionTree with a score of 0.9017857142857142
Best params is : {'criterion': 'gini', 'max_depth': 16, 'max_features': 'auto', 'min_samples_leaf': 4, 'min_samples_split': 10, 'splitter': 'random'}
```
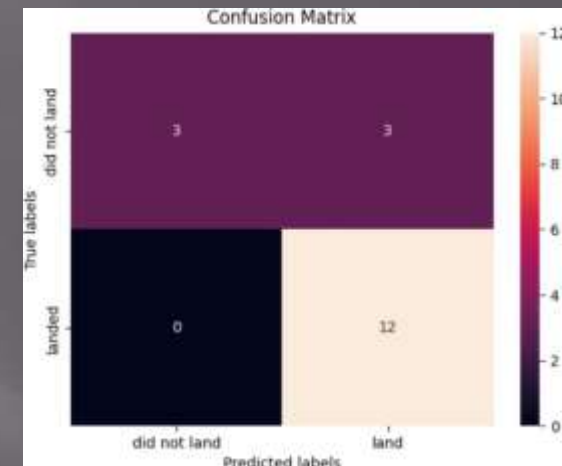
# Confusion Matrices

**Performance Summary**

- A confusion matrix summarizes the performance of a classification algorithm

- All the confusion matrices were identical

- The fact that there are false positives (Type 1 error) is not good

- Confusion Matrix Outputs:

  - 12 True positive
  - 3 True negative
  - **3 False positive**
  - 0 False Negative

- **Precision** = TP / (TP + FP)

  - 12 / 15 = .80

- **Recall** = TP / (TP + FN)

  - 12 / 12 = 1

- **F1 Score** = 2 * (Precision * Recall) / (Precision + Recall)

  - 2 * (.8 * 1) / (.8 + 1) = .89

- **Accuracy** = (TP + TN) / (TP + TN + FP + FN) = .833

# Conclusion

**Research**

- **Model Performance**: The models performed similarly on the test set with the decision tree model slightly outperforming

- **Equator**: Most of the launch sites are near the equator for an additional   natural boost - due to the rotational speed of earth - which helps save the cost of putting in extra fuel and boosters

- **Coast**: All the launch sites are close to the coast

- **Launch Success**: Increases over time

- **KSC LC-39A**: Has the highest success rate among launch sites. Has a 100% success rate for launches less than 5,500 kg

- **Orbits**: ES-L1, GEO, HEO, and SSO have a 100% success rate

- **Payload Mass**: Across all launch sites, the higher the payload mass (kg), the higher the success rate

# Conclusio
## II

**Things to Consider**

- **Dataset**: A larger dataset will help build on the predictive analytics results to help understand if the findings can be generalizable to a larger data set

- **Feature Analysis / PCA**: Additional feature analysis or principal component analysis should be conducted to see if it can help improve accuracy

- **XGBoost**: Is a powerful model which was not utilized in this study. It would be interesting to see if it outperforms the other classification models

2023