# Outline

- Executive Summary

- Introduction

- Methodology

- Results

- Conclusion

# Executive Summary

## Summary of methodologies

- **Collect** data using REST API and web scraping techniques.
- **Wrangle** data to create success/fail outcome variable.
- **Explore** data with data visualization techniques, considering following factors: payload, launch site, flight number and yearly end.
- **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 ranger.
- Build Models to predict landing outcomes using logistic regression, support vector machine(SVM), decision three and k nearest neighbor(KNN).

## Summary of all results

- **Exploratory Data Analysis**: Launch success has improved over time.
- **Visualization**: Most launch sites are near equator.
- **Predictive Analytics**: All models performed similarly on test set.

# Introduction

## Background

SpaceX, a leader in the space industry, strives to make space travel affordable for everyone. Its accomplishments include sending spacecraft to the international space station, launching a satellite constellation that provides internet access and sending manned missions to space. SpaceX can do this because the rocket launches are relatively inexpensive ($62 million per launch) due to its novel reuse of the first stage of its Falcon 9 rocket. Other providers, which are not able to reuse the first stage, cost upwards of $165 million each. By determining if the first stage will land, we can determine the price of the launch. To do this, we can use public data and machine learning models to predict whether SpaceX –or a competing company –can reuse the first stage.

## Explore

- 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

Executive Summary

Data collection methodology: — Collect data using SpaceX REST API and web scraping techniques

Perform data wrangling — Describe how data was processed

Perform exploratory data analysis (EDA) using visualization and SQL

Perform interactive visual analytics using Folium and Plotly Dash

Perform predictive analysis using classification models — To predict landing outcomes using classification models. Tune and evaluate models to find best model and parameters

# Data Collection

- Data collection using SpaceX REST API.

- Data collection using web scraping techniques.

# Data Collection – SpaceX API

**Steps**

- **Request data** from SpaceX API (rocket launch data)

- **Decode response** using .json() and convert to a dataframe using .json_normalize()

- **Request information** about the launches from SpaceX API using custom functions

- **Create dictionary** from the data

- **Create dataframe** from the dictionary

- **Filter dataframe** to contain only Falcon 9 launches

- **Replace missing values** of Payload Mass with calculated .mean()

- **Export data** to csv file

# Data Collection - Scraping

**Steps**

- **Request data** (Falcon 9 launch data) from Wikipedia

- **Create BeautifulSoup object** from HTML response

- **Extract column names** from HTML table header

- **Collect data** from parsing HTML tables

- **Create dictionary** from the data

- **Create dataframe** from the dictionary

- **Export data** to csv file

# Data Wrangling

**1**

Perform EDA and determine data labels

**2**

Calculate and calculate launching outcome

# EDA with Data Visualization

View relationships by using scatter plot . The variables could be useful for machine learning if exists.

Show comparisons among discrete categories with bar charts. Bar charts show the relationships among the categories and measured value.

# EDA with SQL

**Queries**

**Display:**

- Names of unique launch sites

- 5 records where launch site begins with 'CCA'

- Total payload mass carried by boosters launched by NASA (CRS)

- Average payload mass carried by booster version F9 v1.1.

**List:**

- Date of first successful landing on ground pad

- Names of boosters which had success landing on drone ship and have payload mass greater than 4,000 but less than 6,000

- Total number of successful and failed missions

- Names of booster versions which have carried the max payload

- Failed landing outcomes on drone ship, their booster version and launch site for the months in the year 2015

- Count of landing outcomes between 2010-06-04 and 2017-03-20 (desc)

# Build an Interactive Map with Folium

**Markers Indicating Launch Sites**

- Added **blue circle** at **NASA Johnson Space Center's coordinate** with a **popup label** showing its name using its latitude and longitude coordinates

- Added **red circles** at **all launch sites coordinates** with a **popup label** showing its name using its name using its latitude and longitude coordinates

**Colored Markers of Launch Outcomes**

- Added **colored markers** of **successful** (**green**) and **unsuccessful** (**red**) **launches** at each launch site to show which launch sites have high success rates

**Distances Between a Launch Site to Proximities**

- Added **colored lines** to **show distance between** launch site **CCAFS SLC-40 and** its proximity to the **nearest coastline, railway, highway, and city**

# Build a Dashboard with Plotly Dash

**Dropdown List with Launch Sites**

- Allow user to select all launch sites or a certain launch site

**Pie Chart Showing Successful Launches**

- Allow user to see successful and unsuccessful launches as a percent of the total

**Slider of Payload Mass Range**

- Allow user to select payload mass range

**Scatter Chart Showing Payload Mass vs. Success Rate by Booster Version**

- Allow user to see the correlation between Payload and Launch Success

# Predictive Analysis (Classification)

**Charts**

- **Create** NumPy array from the Class column

- **Standardize** the data with StandardScaler. Fit and transform the data.

- **Split** the data using train_test_split

- **Create** a GridSearchCV object with cv=10 for parameter optimization

- **Apply** GridSearchCV on different algorithms: logistic regression (LogisticRegression()), support vector machine (SVC()), decision tree (DecisionTreeClassifier()), K-Nearest Neighbor (KNeighborsClassifier())

- **Calculate** accuracy on the test data using .score() for all models

- **Assess** the confusion matrix for all models

- **Identify** the best model using Jaccard_Score, F1_Score and Accuracy

# Results

- Exploratory data analysis results

- Interactive analytics demo in screenshots

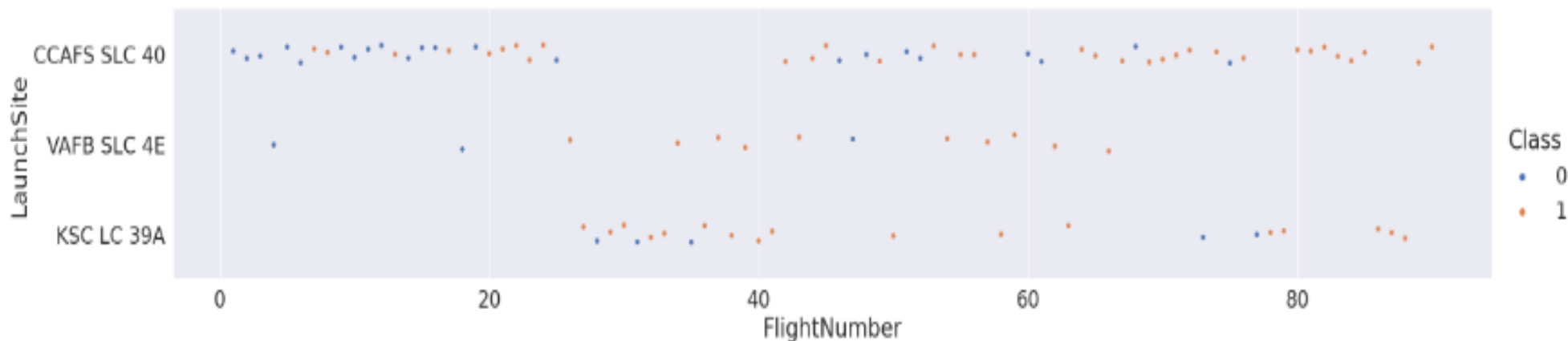- Predictive analysis results

Section 2

# Insights drawn from EDA

# Flight Number vs. Launch Site
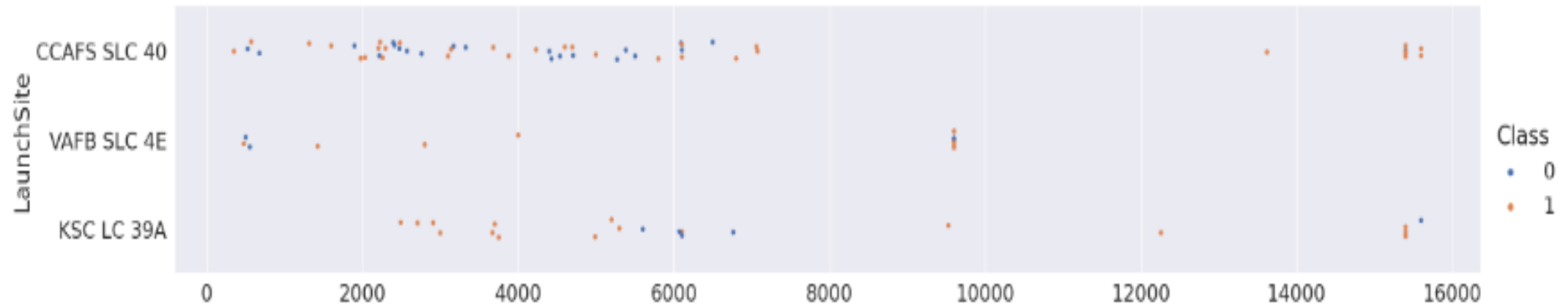
**Exploratory Data Analysis**

- **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
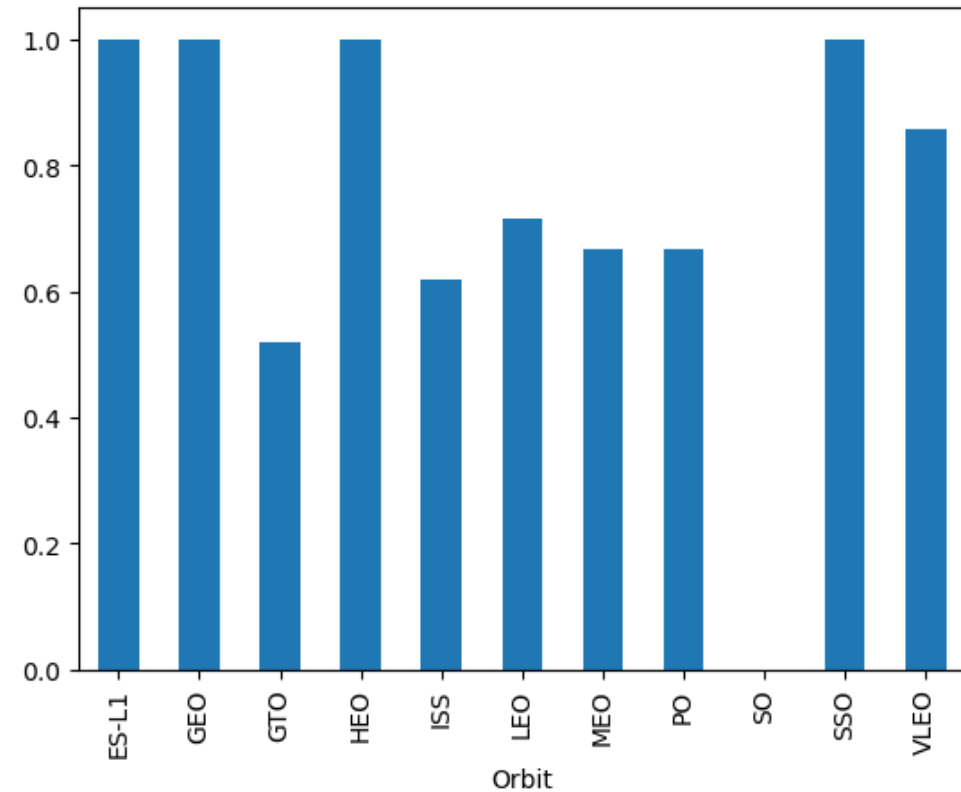
# Payload vs. Launch Site

**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 vs. Orbit Type

**Exploratory Data Analysis**

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

- **50%-80% Success Rate**: GTO, ISS, LEO, MEO, PO
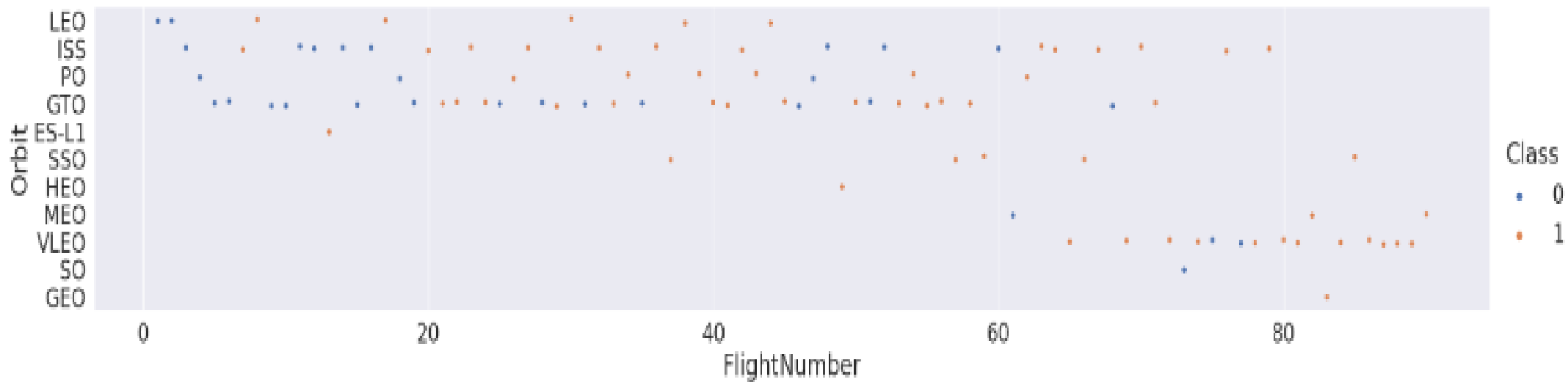
- **0% Success Rate**: SO

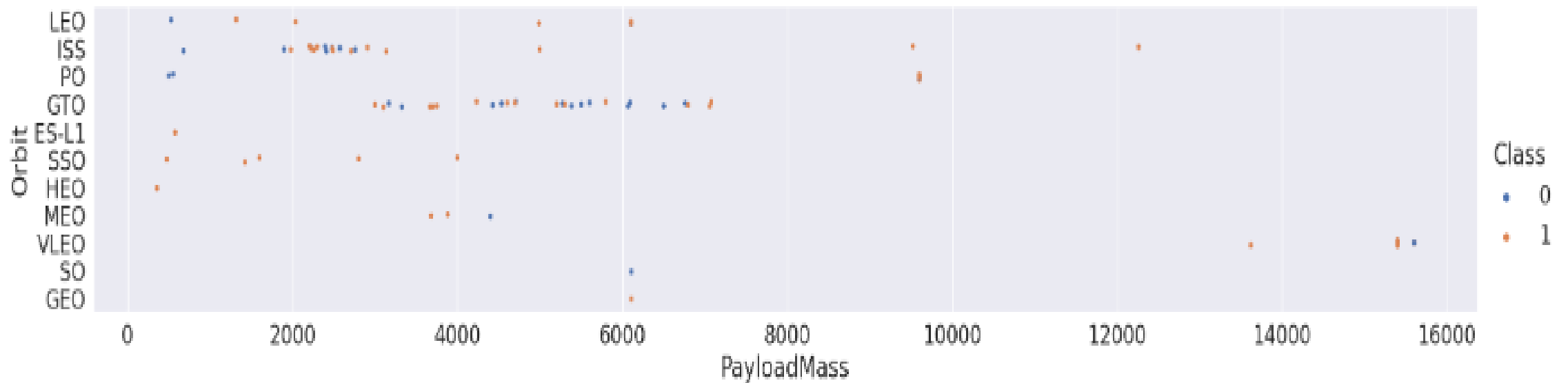# Flight Number vs. Orbit Type

**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

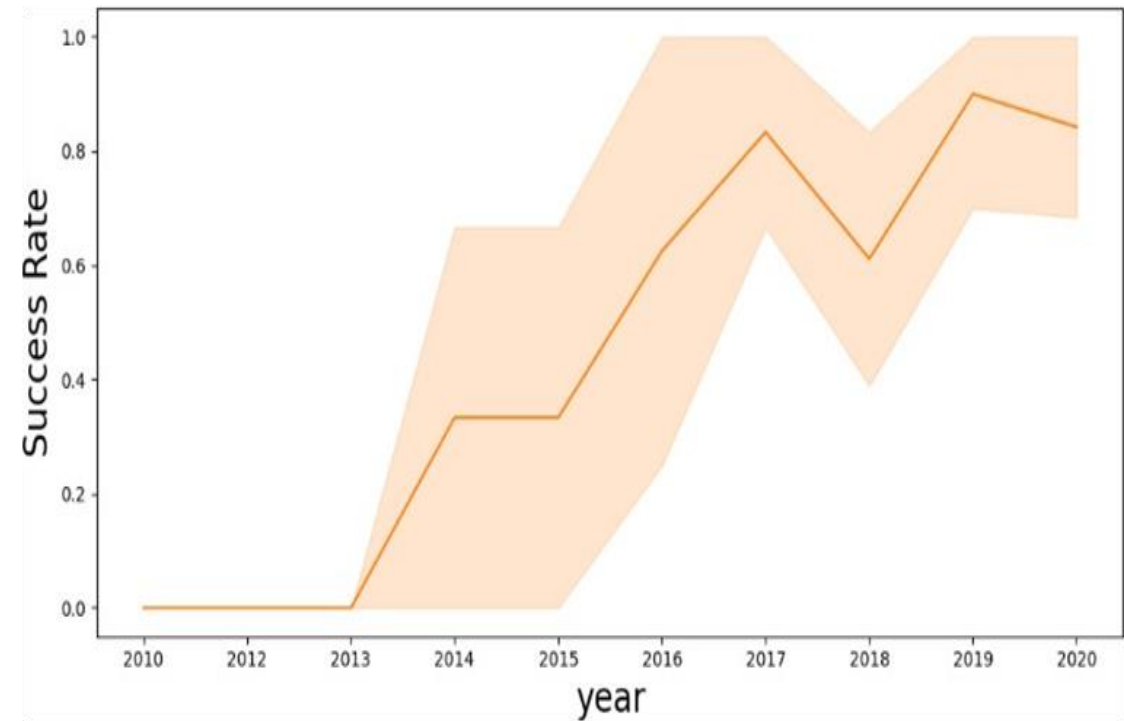# Payload vs. Orbit Type

**Exploratory Data Analysis**

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

- The GTO orbit has mixed success with heavier payloads

# Launch Success Yearly Trend

**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

# All Launch Site Names

## Launch Site Names

- CCAFS LC-40

- CCAFS SLC-40

- KSC LC-39A

- VAFB SLC-4E

```
%sql select distinct Launch_Site from SPACEXTBL
```

* sqlite:///my_data1.db
Done.

**Launch_Site**

CCAFS LC-40

VAFB SLC-4E

KSC LC-39A

CCAFS SLC-40

# Launch Site Names Begin with 'CCA'

## 5 Launch Site Names starts with CCA

- CCAFS SLC-40
- CCAFS LC-40
- CCAFS LC-40
- CCAFS LC-40
- CCAFS LC-40

```
%sql select * from SPACEXTBL where Launch_Site like 'CCA%' limit 5
```
\* sqlite:///my_data1.db
one.

| Date | Time (UTC) | Booster_Version | Launch_Site | Payload | PAYLOAD_MASS__KG_ | Orbit | Customer | Mission_Outcome | Landing_Outcome |
|---|---|---|---|---|---|---|---|---|---|
| 2010-04-06 | 18:45:00 | F9 v1.0 B0003 | CCAFS LC-40 | Dragon Spacecraft Qualification Unit | 0 | LEO | SpaceX | Success | Failure (parachute) |
| 2010-08-12 | 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 | 07:44:00 | F9 v1.0 B0005 | CCAFS LC-40 | Dragon demo flight C2 | 525 | LEO (ISS) | NASA (COTS) | Success | No attempt |
| 2012-08-10 | 00:35:00 | F9 v1.0 B0006 | CCAFS LC-40 | SpaceX CRS-1 | 500 | LEO (ISS) | NASA (CRS) | Success | No attempt |
| 2013-01-03 | 15:10:00 | F9 v1.0 B0007 | CCAFS LC-40 | SpaceX CRS-2 | 677 | LEO (ISS) | NASA (CRS) | Success | No attempt |

# Total Payload Mass

**Total Payload Mass : 45596 kg (total) carried by NASA boosters.**

```
%sql select SUM(PAYLOAD_MASS__KG_) from SPACEXTBL where Customer = 'NASA (CRS)'
```

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

| SUM(PAYLOAD_MASS__KG_) |
|---|
| 45596 |

# Average Payload Mass by F9 v1.1

**Average Payload mass carried by F9 v1.1**: 2928.4

```
%sql select AVG(PAYLOAD_MASS__KG_) from SPACEXTBL where Booster_Version = 'F9 v1.1'
```

 * sqlite:///my_data1.db
Done.

AVG(PAYLOAD_MASS__KG_)

2928.4

# First Successful Ground Landing Date

**The dates of the first successful landing outcome on ground pad**: 2015-12-22

```
%sql select MIN(Date) from SPACEXTBL where Landing_Outcome = 'Success (ground pad)'
```

* sqlite:///my_data1.db
Done.

MIN(Date)

2015-12-22

# Successful Drone Ship Landing with Payload between 4000 and 6000

## Booster Drone Ship Landing

- Booster mass greater than 4,000 but less than 6,000

- JSCAT-14, JSCAT-16, SES-10, SES-11 / EchoStar 105

```
%sql select Booster_Version from SPACEXTBL where Landing_Outcome = 'Success (drone ship)' and PAYLOAD_MASS__KG_ >4000 and P
```

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

| Booster_Version |
| --- |
| F9 FT B1022 |
| F9 FT B1026 |
| F9 FT B1021.2 |
| F9 FT B1031.2 |

# Total Number of Successful and Failure Mission Outcomes

**Total number of successful and failure mission outcomes:** 71

```
%sql select COUNT(*) from SPACEXTBL where Landing_Outcome like 'Success%' or Landing_Outcome like 'Failure%'
```

* sqlite:///my_data1.db
one.

| COUNT(*) |
| --- |
| 71 |

# Boosters Carried Maximum Payload

**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

```
%sql select Booster_Version from SPACEXTBL where PAYLOAD_MASS__KG_ = (select MAX(PAYLOAD_MASS__KG_) from SPACEXTBL)
```

* sqlite:///my_data1.db
Done.

| 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

## In 2015

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

```
%sql select strftime('%m',Date) as "Month",Booster_Version,Launch_Site,Landing_Outcome from SPACEXTBL where strftime('%Y',Da
```

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

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

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

**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.
```

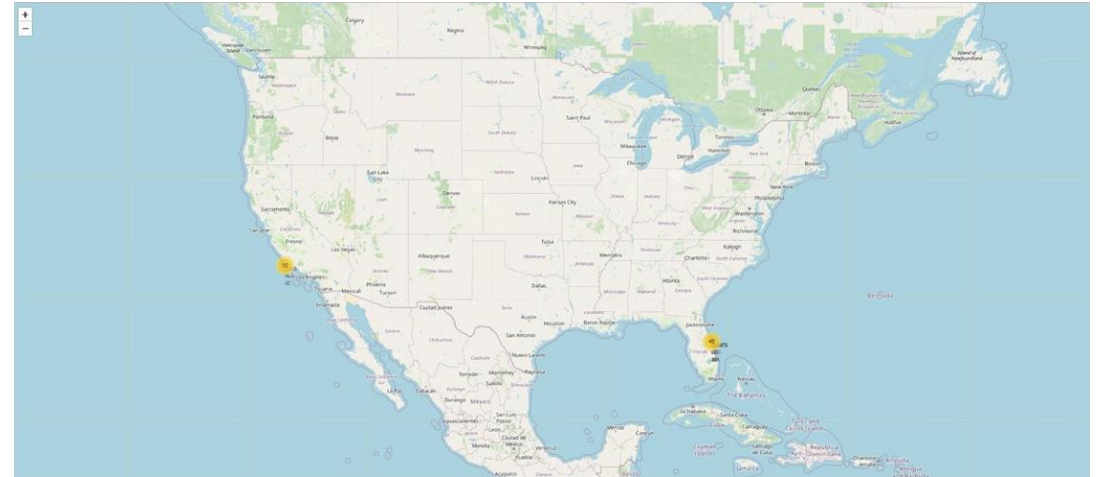| Landing _Outcome | count_outcomes |
|---|---|
| Success | 20 |
| No attempt | 10 |
| Success (drone ship) | 8 |
| Success (ground pad) | 6 |
| Failure (drone ship) | 4 |
| Failure | 3 |
| Controlled (ocean) | 3 |
| Failure (parachute) | 2 |
| No attempt | 1 |

Section 3

# Launch Sites Proximities Analysis

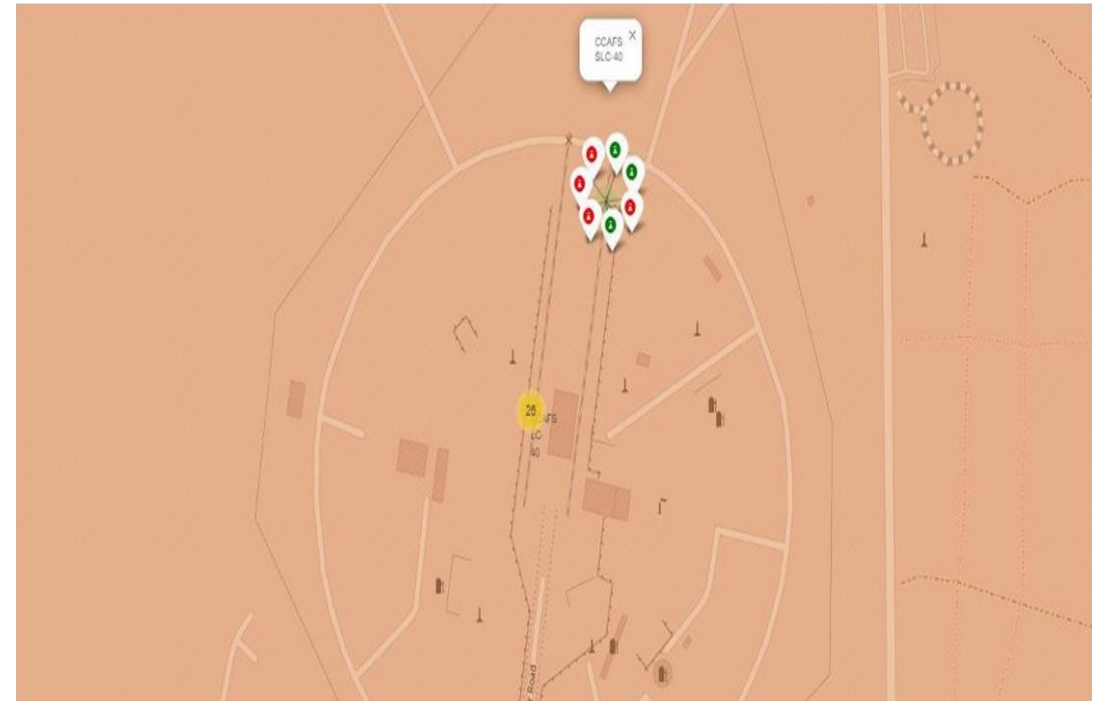# Launch Sites

**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.
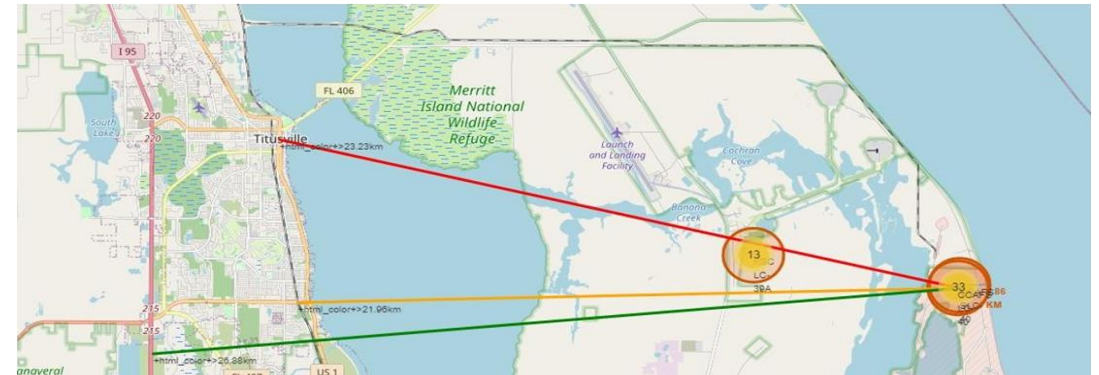
# Launch Outcomes

**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

# Build a Dashboard with Plotly Dash

# Launch Success by Site

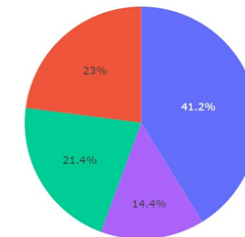**Success as Percent of Total**

- **KSC LC-39A** has the **most successful launches** amongst launch sites (**41.2%**)



**SpaceX Launch Records Dashboard**

All Sites

Total Success Launches by Site

- KSC LC-39A
- CCAFS SLC-40
- VAFB SLC-4E
- CCAFS LC-40
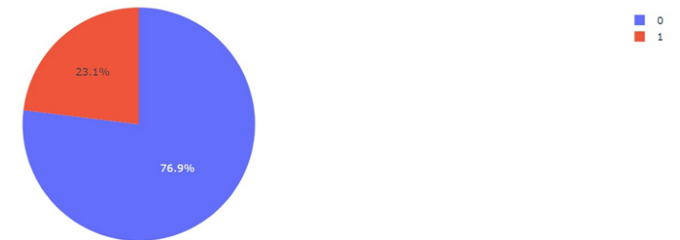
41.2%
23%
21.4%
14.4%

# 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
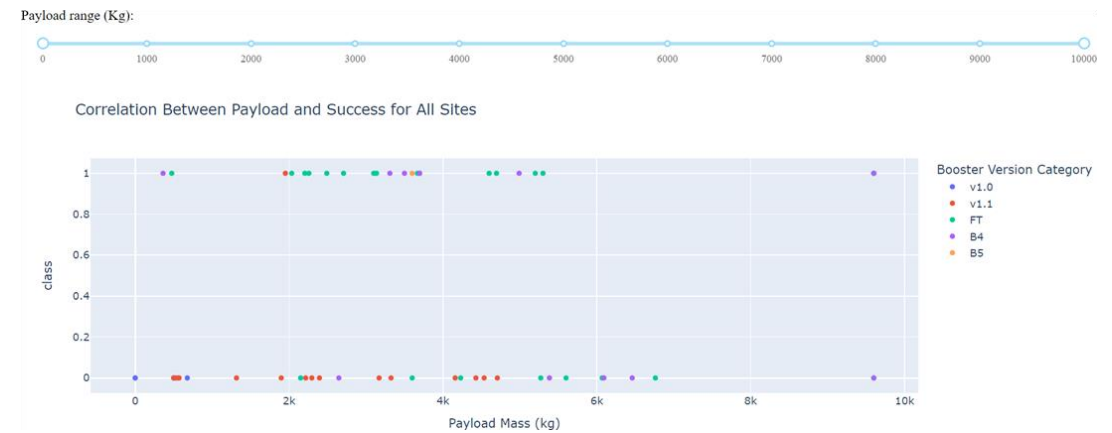


SpaceX Launch Records Dashboard

KSC LC-39A

Total Success Launches for Site KSC LC-39A

23.1%

76.9%

0
1

# 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 Analysis (Classification)

# Classification Accuracy

**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

```
3]:   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])
```
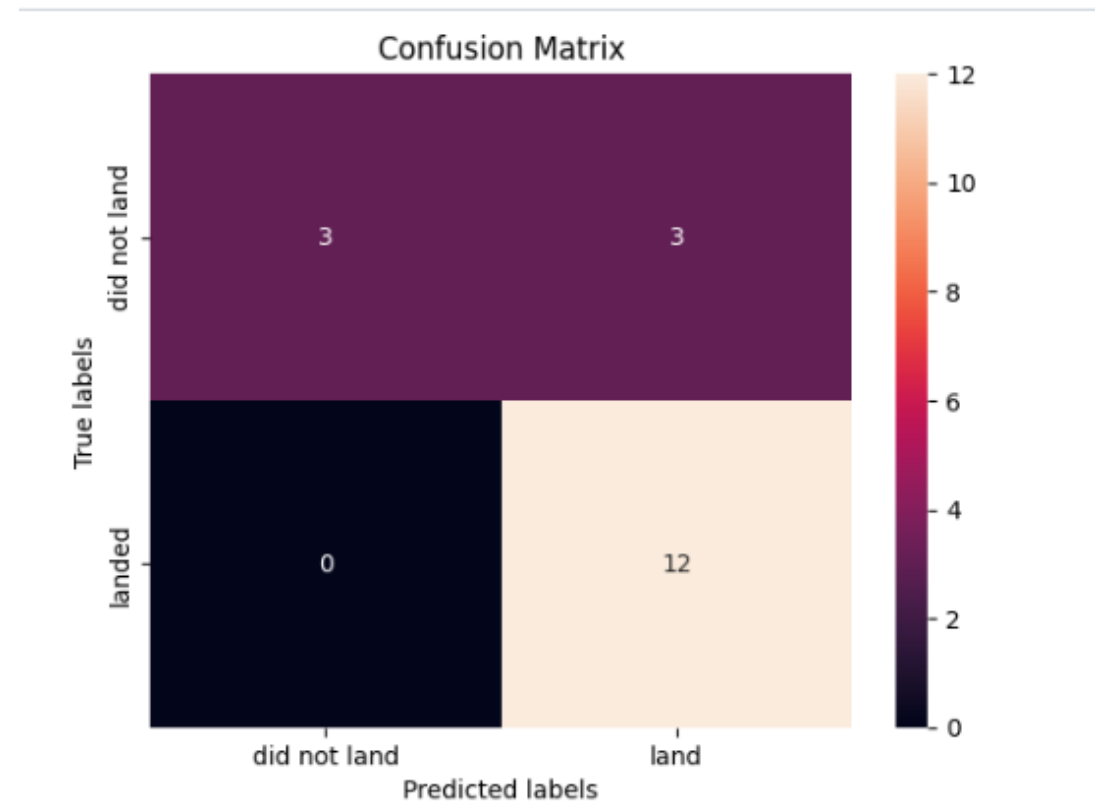
Best model is DecisionTree with a score of 0.8767857142857143

# Confusion Matrix

**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



Confusion Matrix

# Conclusions

- **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

Thank you!