IBM Developer
SKILLS NETWORK

# Winning Space Race with Data Science

GAURAV M. DABADE
01 / 08 / 24

# Outline

- Executive Summary

- Introduction

- Methodology

- Results

- Conclusion

- Appendix

# Executive Summary

- Data Collection – Space X API & SPACE X WIKI
- Four ML Model are used -
    - Logistic Regression
    - SVM
    - Decision Tree Classifier
    - KNN
- Exported Data from SQL
- Project uses Folium Maps and Visualizations Methods
- The Accuracy Rate of Successful Landings – 83.33%

# Introduction

- Aim – SPACE Y is about to compete with SPPACE X
- Service – Commercial Space Carrier
- Problem – To tackle the landings failure using Data Science
- Solution -
    - Study and Analyze SPACE X DATA
    - Calculate the findings and Deploy the Solutions at SPACE Y.

Section 1

# Methodology

# Methodology

## Executive Summary

- Data collection methodology:

  - Data collected by using SPACE X API and Web Scraping

- Perform data wrangling

  - By classifying successful landings and unsuccessful landings

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

- Perform interactive visual analytics using Folium and Plotly Dash

- Perform predictive analysis using classification models

  - Tuned Models using GridSearchCV
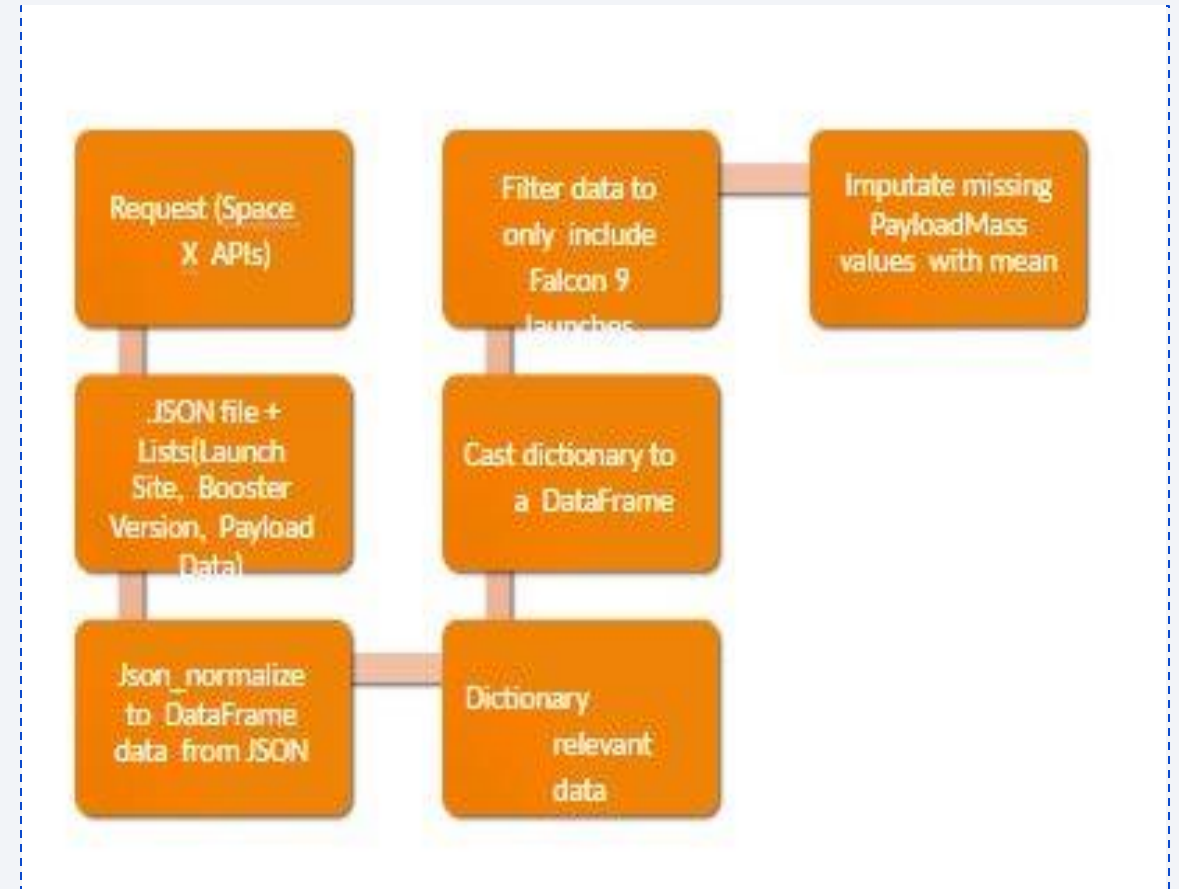
# Data Collection

- Data collection process involved a combination of API requests from Space X public API and web scraping data from a table in Space X's Wikipedia entry.

- The next slide will show the flowchart of data collection from API and the one after will show the flowchart of data collection from web scraping.

- Space X API Data Columns:

- FlightNumber, Date, BoosterVersion, PayloadMass, Orbit, LaunchSite, Outcome, Flights, GridFins,

- Reused, Legs, LandingPad, Block, ReusedCount, Serial, Longitude, Latitude

- Wikipedia Webscrape Data Columns:

- Flight No., Launch site, Payload, PayloadMass, Orbit, Customer, Launch outcome, Version Booster, Booster landing, Date, Time

# Data Collection – SpaceX API

SPACE X DATA COLLECTION

# Data Collection - Scraping

- [SPACE X DATA WEB SCRAPING](#)

# Data Wrangling

- Create a training label with landing outcomes where successful = 1 & failure = 0.

- Outcome column has two components: 'Mission Outcome' 'Landing Location'

- New training label column 'class' with a value of 1 if 'Mission Outcome' is True and 0 otherwise.  Value Mapping:

- True ASDS, True RTLS, & True Ocean – set to -> 1

- None None, False ASDS, None ASDS, False Ocean, False RTLS – set to -> 0

DATA WRANGLING

# EDA with Data Visualization

- Exploratory Data Analysis performed on variables Flight Number, Payload Mass, Launch Site,  Orbit, Class and Year.

- Plots Used:

- Flight Number vs. Payload Mass, Flight Number vs. Launch Site, Payload Mass vs. Launch Site,  Orbit vs. Success Rate, Flight Number vs. Orbit, Payload vs Orbit, and Success Yearly Trend

- Scatter plots, line charts, and bar plots were used to compare relationships between variables to

- decide if a relationship exists so that they could be used in training the machine learning model

- EDA DATA VISUALIZATION

# EDA with SQL

- Loaded data set into IBM DB2 Database.

- Queried using SQL Python integration.

- Queries were made to get a better understanding of the dataset.

- Queried information about launch site names, mission outcomes, various pay load sizes of customers and booster versions, and landing outcomes

[EDA WITH SQL](#)

# Build an Interactive Map with Folium

- Folium maps mark Launch Sites, successful and unsuccessful landings, and a proximity example  to key locations: Railway, Highway, Coast, and City.

- This allows us to understand why launch sites may be located where they are. Also visualizes  successful landings relative to location.

- [FOLLIUM MAPS](#)

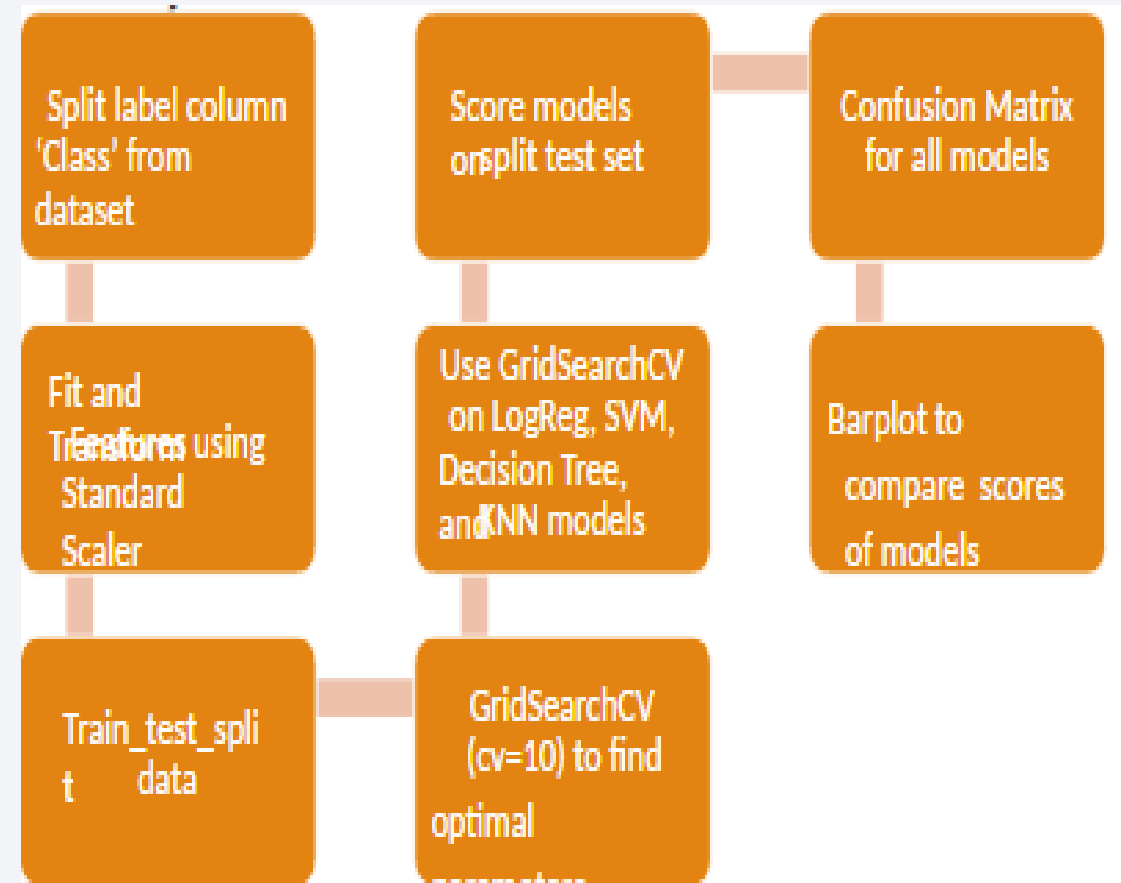# Build a Dashboard with Plotly Dash

- Dashboard includes a pie chart and a scatter plot.
- Pie chart can be selected to show distribution of successful landings across all launch sites and  can be selected to show individual launch site success rates.
- Scatter plot takes two inputs: All sites or individual site and payload mass on a slider between 0  and 10000 kg.
- The pie chart is used to visualize launch site success rate.
- The scatter plot can help us see how success varies across launch sites, payload mass, and
- booster version category.


- [PLOTLY AND DASH -  DASHBOARD](#)

# Predictive Analysis (Classification)

## ML CLASSIFICATION AND PREDICTION

# Results

- Exploratory data analysis results

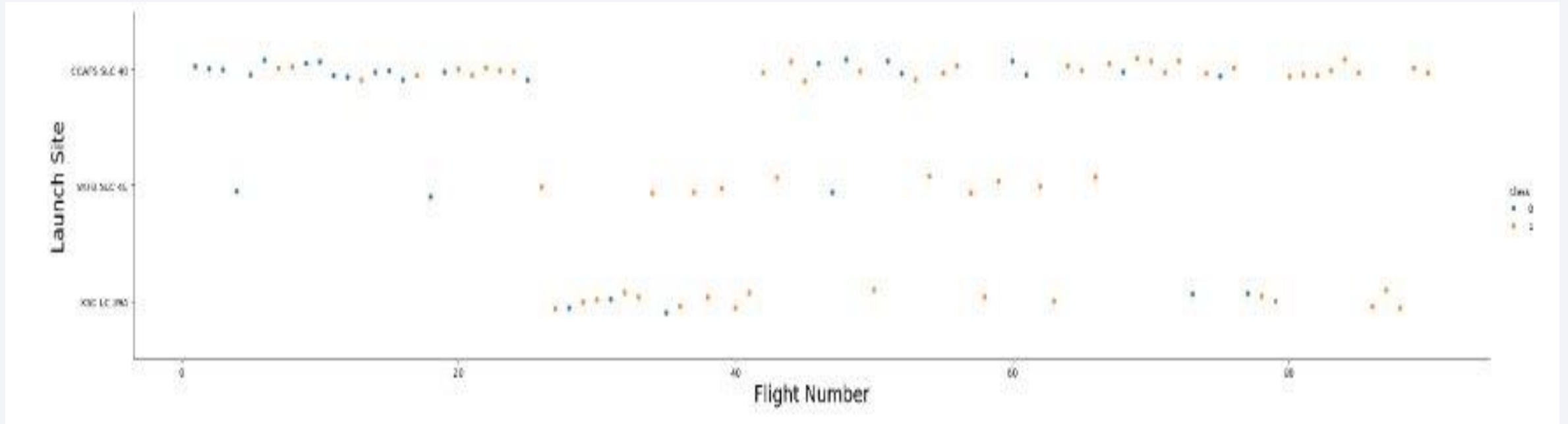- Interactive analytics demo in screenshots

- Predictive analysis results

Section 2

# Insights drawn from EDA

# Flight Number vs. Launch Site



Graphic suggests an increase in success rate over time (indicated in Flight Number).  Likely a big breakthrough around flight 20 which significantly increased success rate.  CCAFS appears to be the main launch site as it has the most volume.
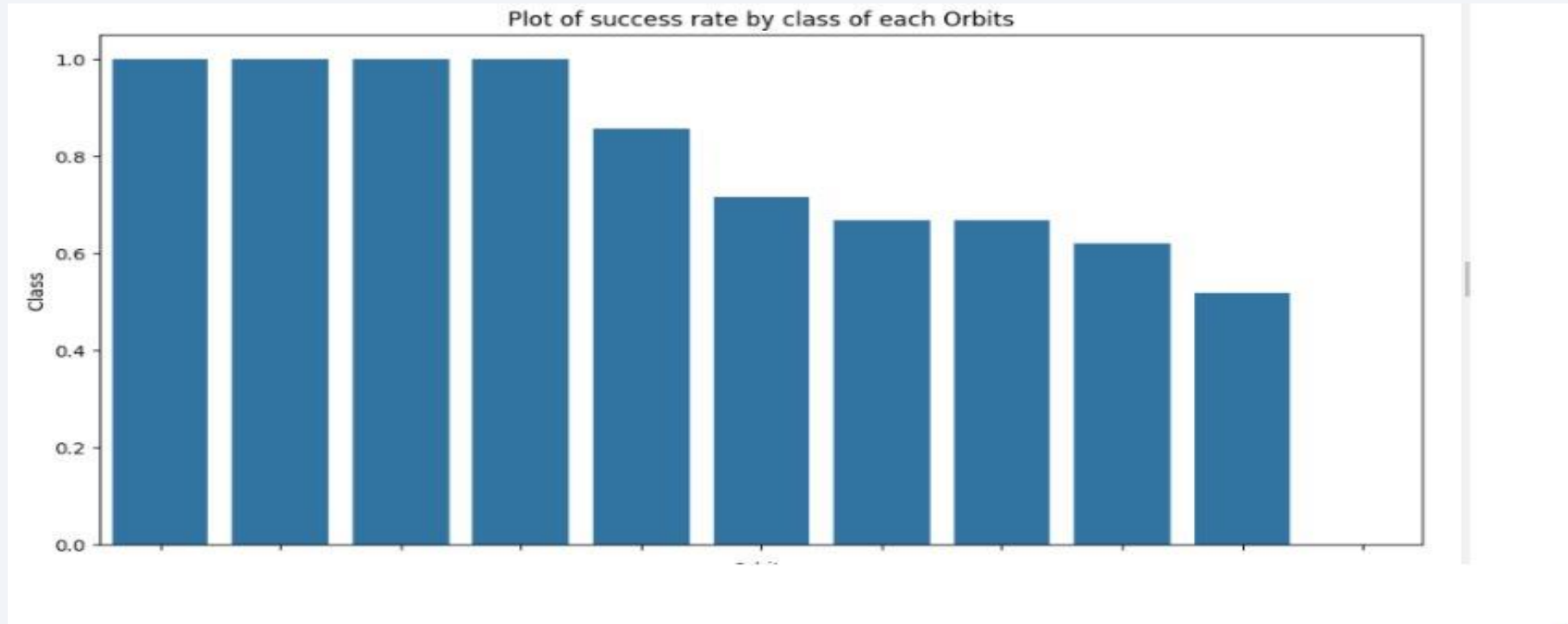
# Payload vs. Launch Site



Payload mass appears to fall mostly between 0-6000 kg.  Different launch sites also seem to use different payload mass.

# Success Rate vs. Orbit Type
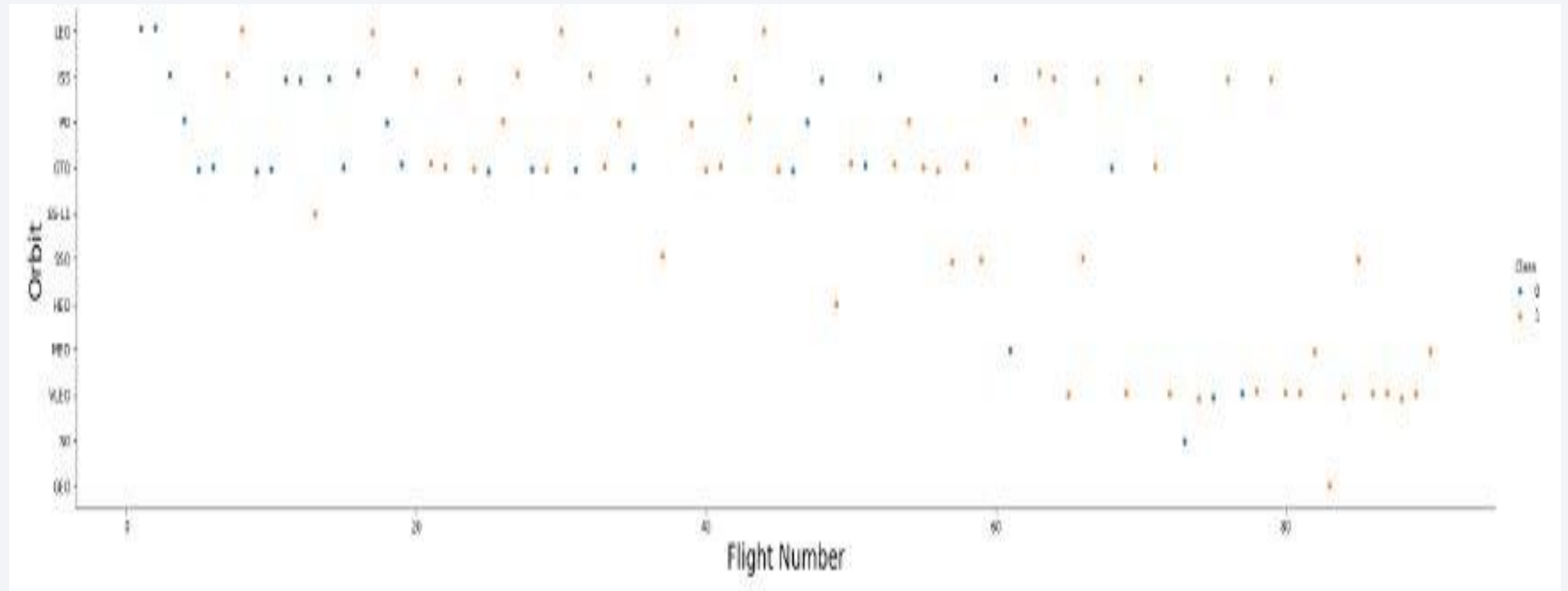


Plot of success rate by class of each Orbits

ES-L1 (1), GEO (1), HEO (1) have 100% success rate (sample sizes in parenthesis)  SSO (5) has 100% success rate
VLEO (14) has decent success rate and attempts
SO (1) has 0% success rate
GTO (27) has the around 50% success rate but largest sample
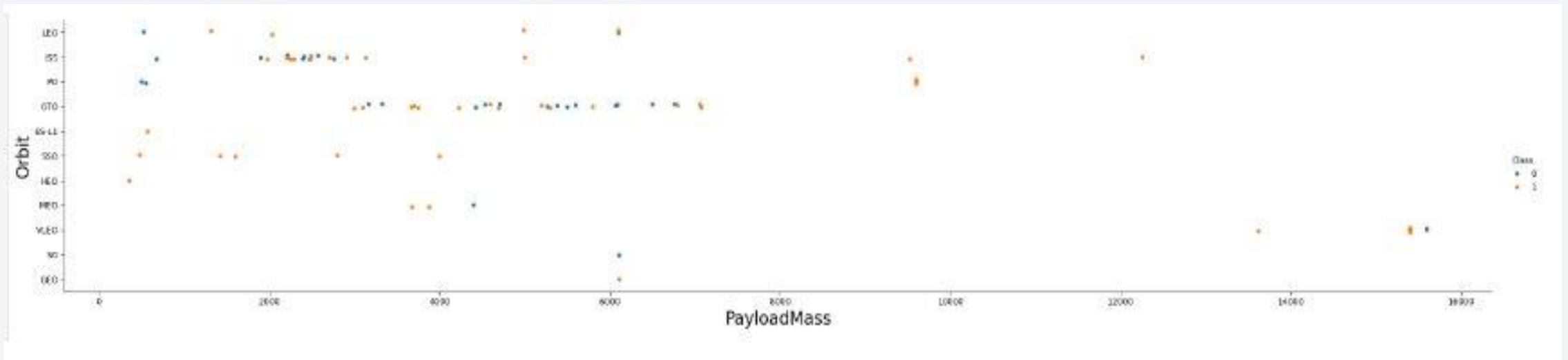
# Flight Number vs. Orbit Type



Launch Orbit preferences changed over Flight Number.  Launch Outcome seems to correlate with this preference.
SpaceX started with LEO orbits which saw moderate success LEO and returned to VLEO in recent launches  SpaceX appears to perform better in lower orbits or Sun-synchronous orbits
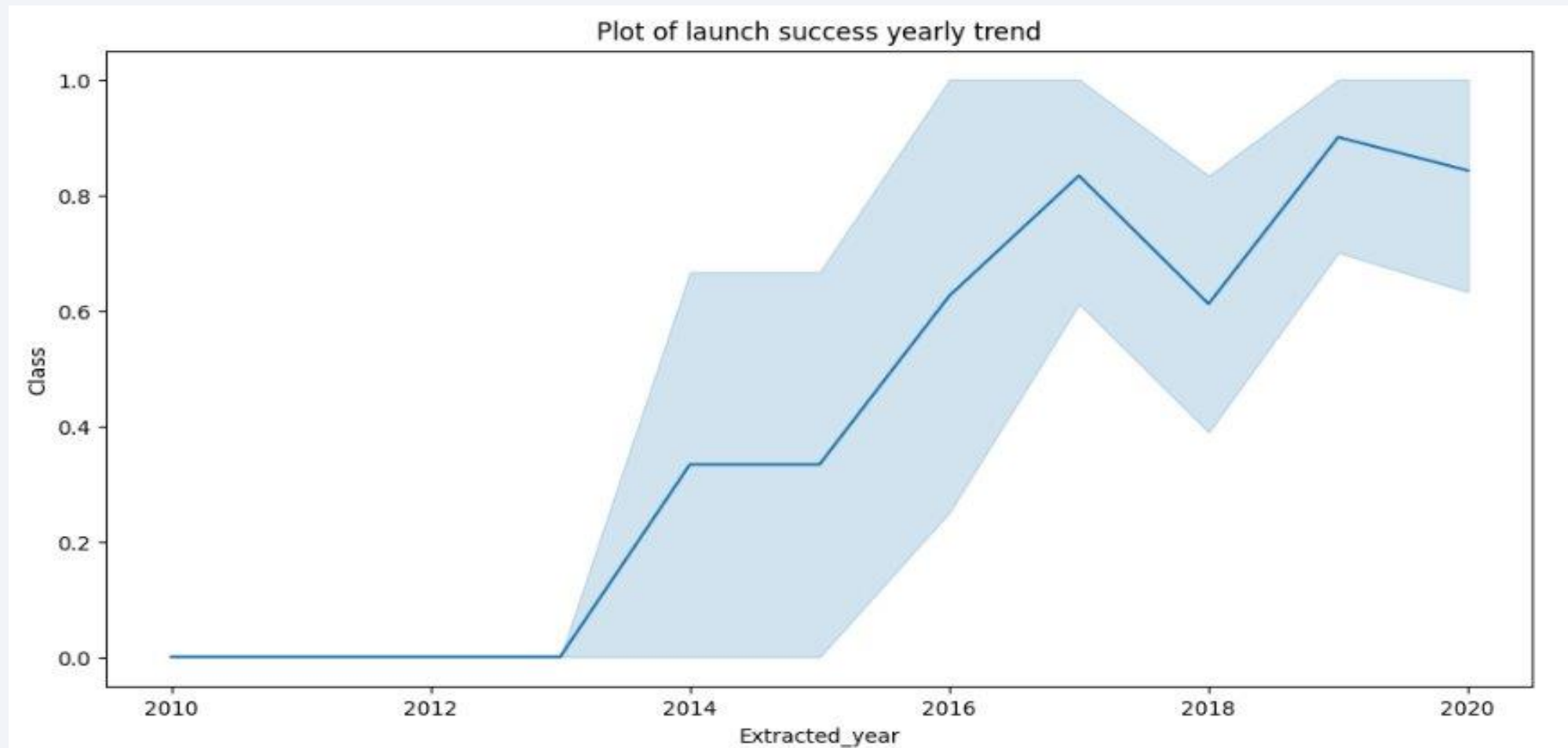
# Payload vs. Orbit Type



Payload mass seems to correlate with orbit
LEO and SSO seem to have relatively low payload mass
The other most successful orbit VLEO only has payload mass values in the higher end of the range

# Launch Success Yearly Trend



Plot of launch success yearly trend

# All Launch Site Names

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

- There are four launch sites .

# Launch Site Names Begin with 'CCA'

| Booster_Version | Launch_Site |
| --- | --- |
| F9 v1.0 B0003 | CCAFS LC-40 |
| F9 v1.0 B0004 | CCAFS LC-40 |
| F9 v1.0 B0005 | CCAFS LC-40 |
| F9 v1.0 B0006 | CCAFS LC-40 |
| F9 v1.0 B0007 | CCAFS LC-40 |

First five entries  in database with  Launch Site name  beginning with  CCA.

# Total Payload Mass

**sum**

**45596**

This query sums the total payload  mass in kg where NASA was the  customer.
CRS stands for Commercial  Resupply Services which indicates  that these payloads were sent to  the International Space Station  (ISS).

# Average Payload Mass by F9 v1.1

**Average**

2534.6666666666665

This query calculates the average payload mass or launches which used booster version F9 v1.1 Average payload mass of F9 1.1 is on the low end of our payload mass range

# First Successful Ground Landing Date

| Date |
| --- |
| 2010-06-04 |

This query returns the first  successful ground pad landing  date.

# Successful Drone Ship Landing with Payload between 4000 and 6000

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

This query returns the four  booster versions that had  successful drone ship landings  and a payload mass between  4000 and 6000 non-inclusively.

# Total Number of Successful and Failure Mission Outcomes

| Mission_Outcome | Count |
|---|---|
| Failure (in flight) | 1 |
| Success | 98 |
| Success | 1 |
| Success (payload status unclear) | 1 |

This query returns a count of each mission outcome.
SpaceX appears to achieve its mission outcome nearly 98% of the time.

# Boosters Carried Maximum 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 |

This likely indicates payload mass
correlates  with the booster version that is used.

# 2015 Launch Records

| MONTH | landing__outcome | booster_version | payload_mass__kg_ | launch_site |
|---|---|---|---|---|
| January | Failure (drone ship) | F9 v1.1 B1012 | 2395 | CCAFS LC-40 |
| April | Failure (drone ship) | F9 v1.1 B1015 | 1898 | CCAFS LC-40 |

This query returns the Month, Landing  Outcome, Booster Version, Payload  Mass (kg), and Launch site of 2015  launches where stage 1 failed to land  on a drone ship. There were two such occurrences.

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

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

- This query returns a list of successful landings  and between 2010-06-04 and 2017-03-20  inclusively.
- There are 10 no attempts.
- Total success is about 10 rank.

Section 3

# Launch Sites Proximities Analysis

# <Folium Map Screenshot 1>



The left map shows all launch sites relative US map. The right map shows the two Florida launch  sites since they are very close to each other. All launch sites are near the ocean.

# <Folium Map Screenshot 2>



Clusters on Folium map can be clicked on to display each successful
landing (green icon) and failed
landing (red icon). In this example VAFB SLC-4E shows 4 successful
landings and 6 failed landings.

# <Folium Map Screenshot 3>



Launch sites are very close to railways for large part and supply transportation. Launch sites are close to highways for human and supply transport. Launch sites are also close to coasts and relatively far from cities so that launch failures can land in the sea to avoid rockets falling on densely populated areas.
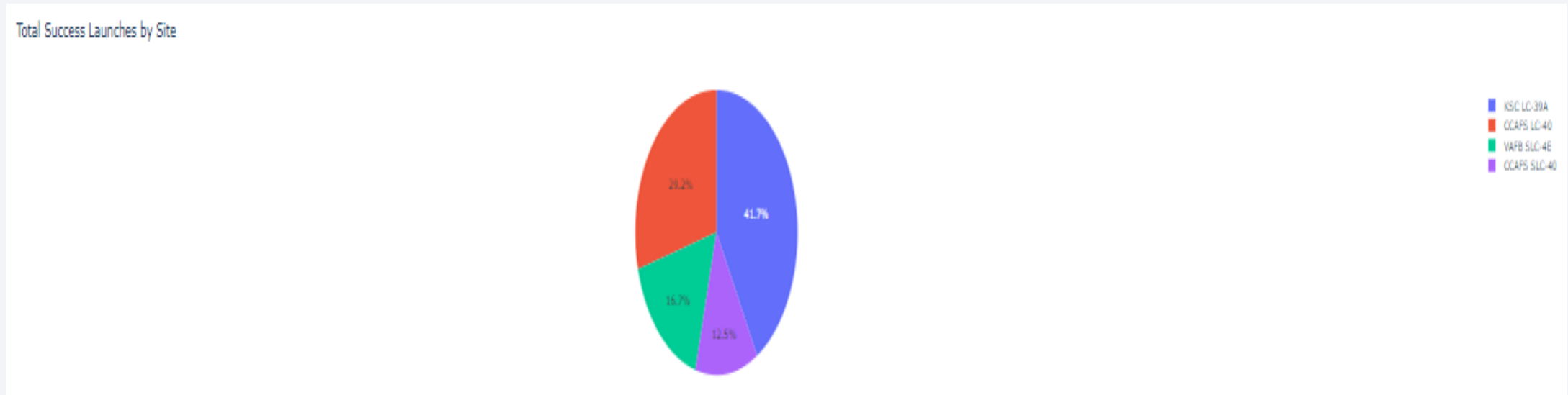
# Build a Dashboard with Plotly Dash

# <Dashboard Screenshot 1>



Total Success Launches by Site

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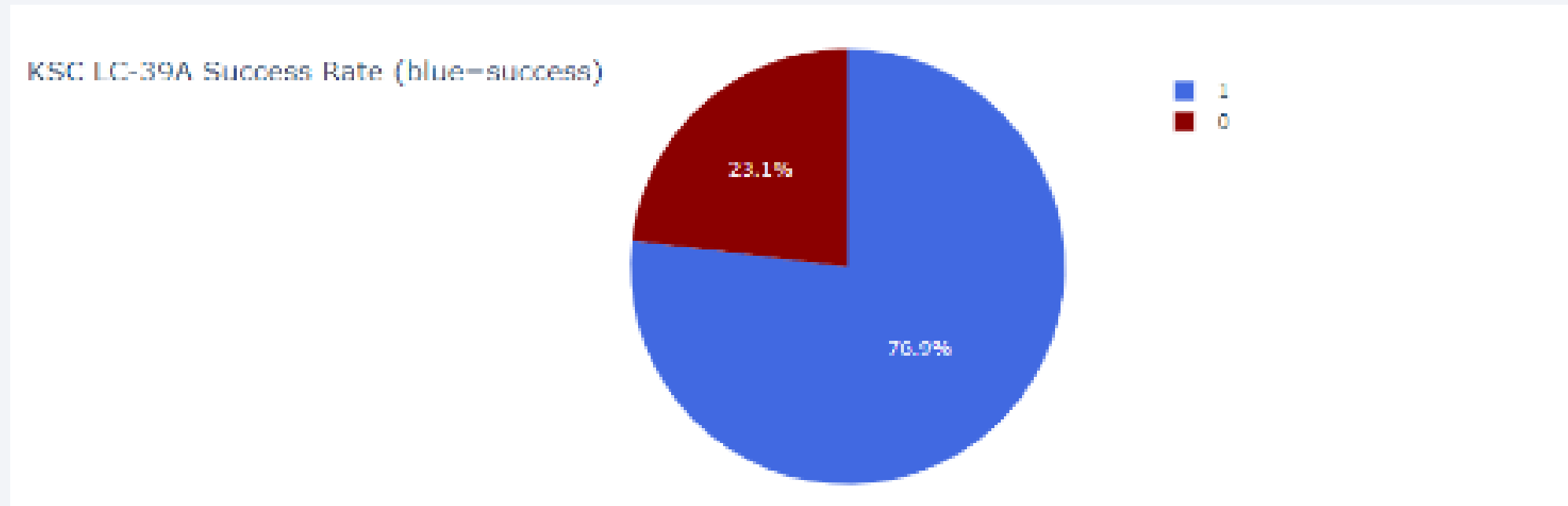CCAFS LC-40
VAFB SLC-4E
CCAFS SLC-40

41.7%
29.2%
16.7%
12.5%

This is the distribution of successful landings across all launch sites. CCAFS LC-40 is the old name of  CCAFS SLC-40 so CCAFS and KSC have the same amount of successful landings, but a majority of the  successful landings where performed before the name change. VAFB has the smallest share of successful  landings. This may be due to smaller sample and increase in difficulty of launching in the west coast.
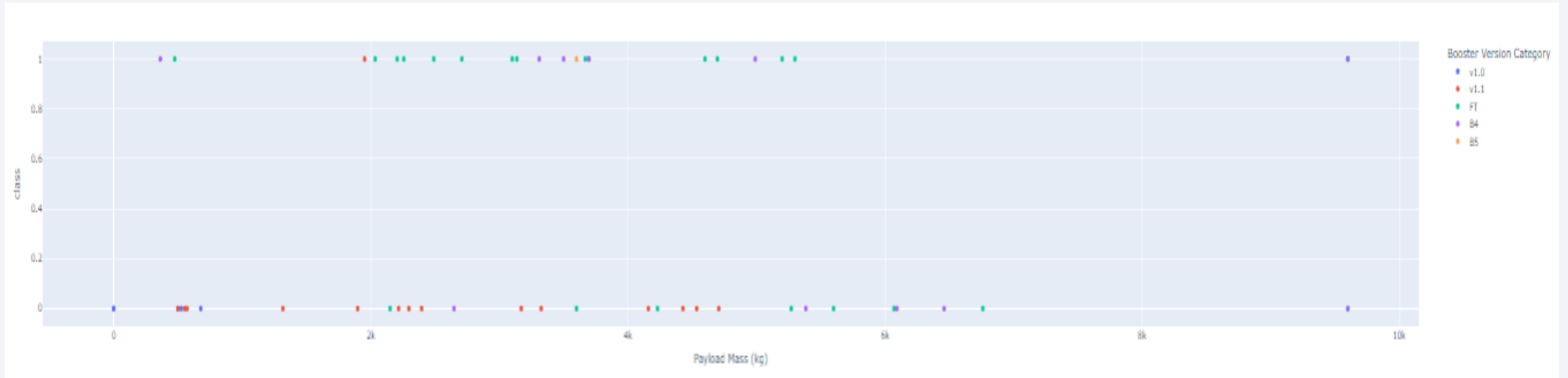
39

# <Dashboard Screenshot 2>



KSC LC-39A has the highest success rate with 10 successful landings and 3 failed landings.

# <Dashboard Screenshot 3>



Plotly dashboard has a Payload range selector. However, this is set from 0-10000 instead of the  max Payload of 15600. Class indicates 1 for successful landing and 0 for failure. Scatter plot also  accounts for booster version category in color and number of launches in point size. In this  particular range of 0-6000, interestingly there are two failed landings with payloads of zero kg.
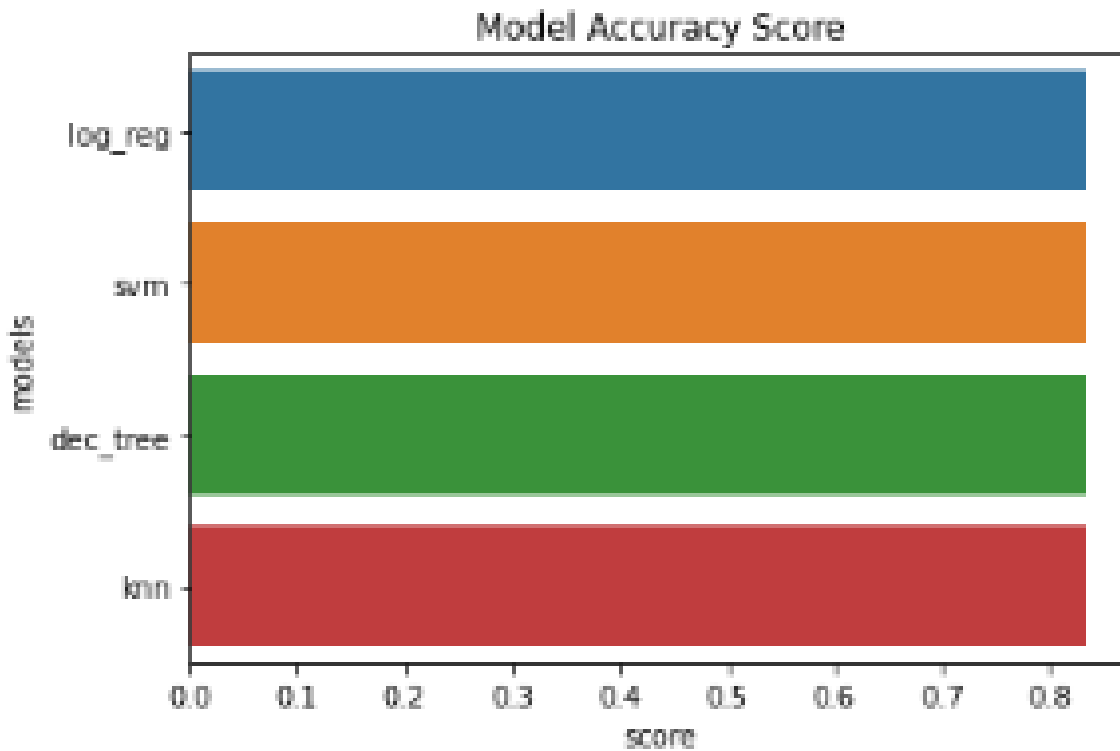
Section 5

# Predictive Analysis (Classification)

# Classification Accuracy
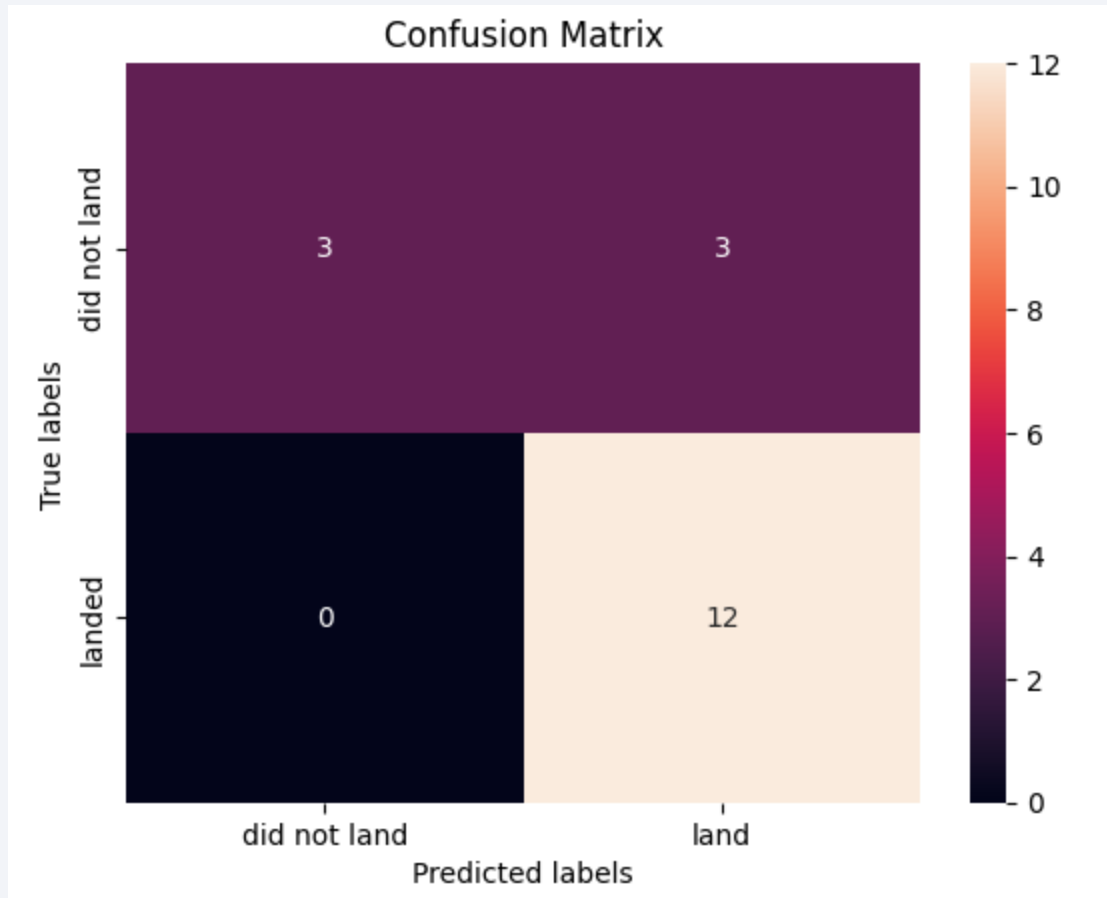


Model Accuracy Score

All models had virtually the same accuracy on the test set at 83.33% accuracy.  It should be noted that test size is small at only sample size of 18.

This can cause large variance in accuracy results, such as those in Decision Tree Classifier model in repeated runs.

We likely need more data to determine the best model.

# Confusion Matrix



Confusion Matrix

Since all models performed the same for the test set, the confusion matrix is the same across all models. The models predicted 12 successful landings when the true label was successful landing.

The models predicted 3 unsuccessful landings when the true label was unsuccessful landing. The models predicted 3 successful landings when the true label was unsuccessful landings (false positives). Our models over predict successful landings.

# Conclusions

- Allon Mask can save 100 million USD.

- He can use this model for better accuracy.

- We can improve accuracy with more data.

- All the data is collected by using SPACE X API AND Web scraping

# Appendix

```python
parameters ={'C':[0.01,0.1,1],
             'penalty':['l2'],
             'solver':['lbfgs']}
```

```python
parameters ={"C":[0.01,0.1,1],'penalty':['l2'], 'solver':['lbfgs']}# l1 lasso l2 ridge
lr=LogisticRegression()
logreg_cv = GridSearchCV(estimator=lr, cv=10, param_grid=parameters).fit(X_train, Y_train)
```

```python
parameters = {'criterion': ['gini', 'entropy'],
              'splitter': ['best', 'random'],
              'max_depth': [2*n for n in range(1,10)],
              'max_features': ['auto', 'sqrt'],
              'min_samples_leaf': [1, 2, 4],
              'min_samples_split': [2, 5, 10]}

tree = DecisionTreeClassifier()
tree_cv = GridSearchCV(estimator=tree, cv=10, param_grid=parameters).fit(X_train, Y_train)
```

```sql
%sql select sum(payload_mass__kg_) as sum from SPACEXTABLE where customer like 'NASA (CRS)'
```

Thank you!