



# Using Data Science to win a space race

Yves-Daniel Toti

11th June 2022

# OUTLINE

---



- Executive Summary
- Introduction
- Methodology
- Results
  - Visualization – Charts
  - Dashboard
- Discussion
  - Findings & Implications
- Conclusion
- Appendix

# EXECUTIVE SUMMARY

---



- Summary of methodologies
  - Data Collection through API
  - Data Collection with Web Scraping
  - Data Wrangling
  - Exploratory Data Analysis with SQL
  - Exploratory Data Analysis with Data Visualization
  - Interactive Visual Analytics with Folium
  - Machine Learning Prediction
- Summary of all results
  - Exploratory Data Analysis result
  - Interactive analytics in screenshots
  - Predictive Analytics results

# INTRODUCTION

---

- Project background and context

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 we are trying to find answers for :
  - What factors determine if the rocket will land successfully?
  - The interaction amongst various features that determine the success rate of a successful -landing.
  - What operating conditions need to be in place to ensure a successful landing program.

# METHODOLOGY

---

## Executive Summary

- Data collection methodology:
  - Data was collected using SpaceX API and web scraping from Wikipedia.
- Perform data wrangling
  - One-hot encoding was applied to categorical features
- Perform exploratory data analysis (EDA) using visualization and SQL
- Perform interactive visual analytics using Folium and Plotly Dash
- Perform predictive analysis using classification models
  - How to build, tune, evaluate classification models

# 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 – SpaceX 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  
: <https://github.com/dannytoti/DataScienceCapstone/blob/a260cd1275a91a3119d1510a615aa97bad932434/jupyter-labs-spacex-data-collection-api.ipynb>

## 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/dannytoti/DataScienceCapstone/blob/32353af519c2868cb98d99c2615a54fa84b5ebb1/jupyter-labs-webscraping.ipynb>



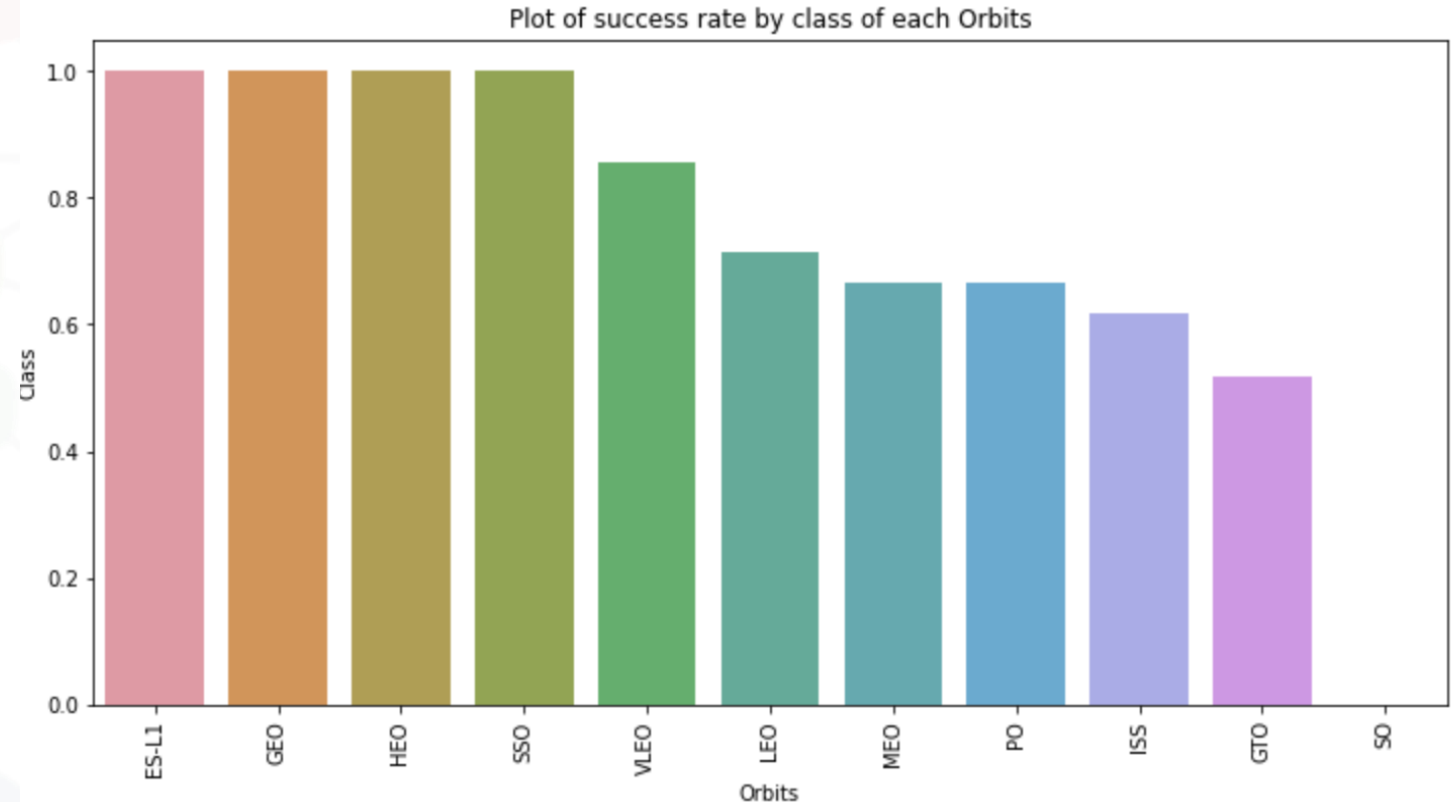
# Data Wrangling

---

- We performed exploratory data analysis and determined the training labels.
- We calculated the number of launches at each site, and the number and occurrence of each orbits
- We created landing outcome label from outcome column and exported the results to csv.

# EDA with Data 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.



# EDA with SQL

---

- We loaded the SpaceX dataset into a **DB2** database
- 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/dannytoti/DataScienceCapstone/blob/32353af519c2868cb98d99c2615a54fa84b5ebb1/jupyter-labs-eda-sql-coursera%20\(1\).ipynb](https://github.com/dannytoti/DataScienceCapstone/blob/32353af519c2868cb98d99c2615a54fa84b5ebb1/jupyter-labs-eda-sql-coursera%20(1).ipynb)

# Build an Interactive 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.
- Using the color-labeled marker clusters, we identified which launch sites have relatively high success rate.

# Build a 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/dannytoti/DataScienceCapstone/blob/5ac637a4c249c842760a6b75c25e04ddbad5995b/spaceXdash.py>

# Predictive Analysis (Classification)

---

- 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/dannytoti/DataScienceCapstone/blob/5ac637a4c249c842760a6b75c25e04ddbad5995b/SpaceX\\_Machine%20Learning%20Prediction\\_Part\\_5.ipynb](https://github.com/dannytoti/DataScienceCapstone/blob/5ac637a4c249c842760a6b75c25e04ddbad5995b/SpaceX_Machine%20Learning%20Prediction_Part_5.ipynb)

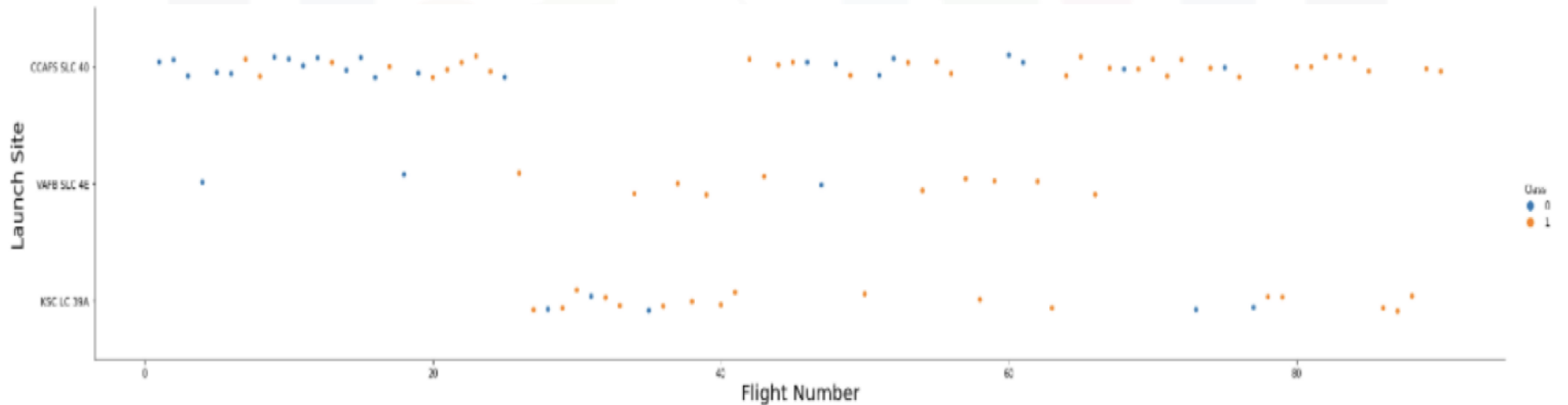
# Result

---

- Exploratory data analysis results
- Interactive analytics demo in screenshots
- Predictive analysis results

# Flight Number vs. Launch Site

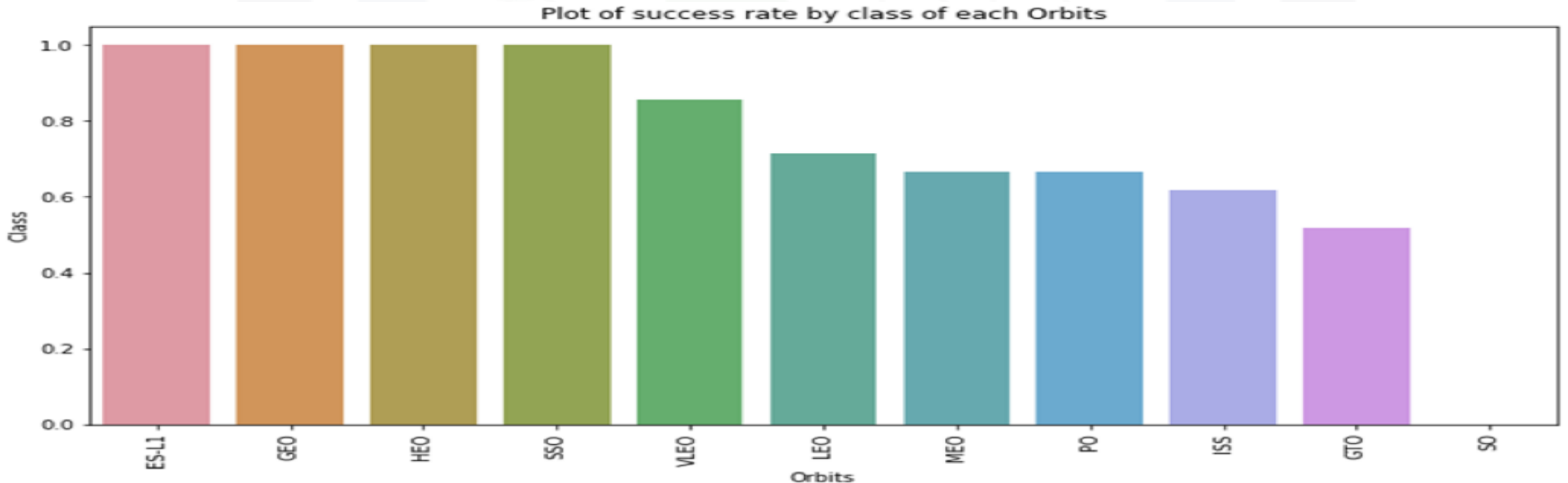
- From the plot, we found that the larger the flight amount at a launch site, the greater the success rate at a launch site.





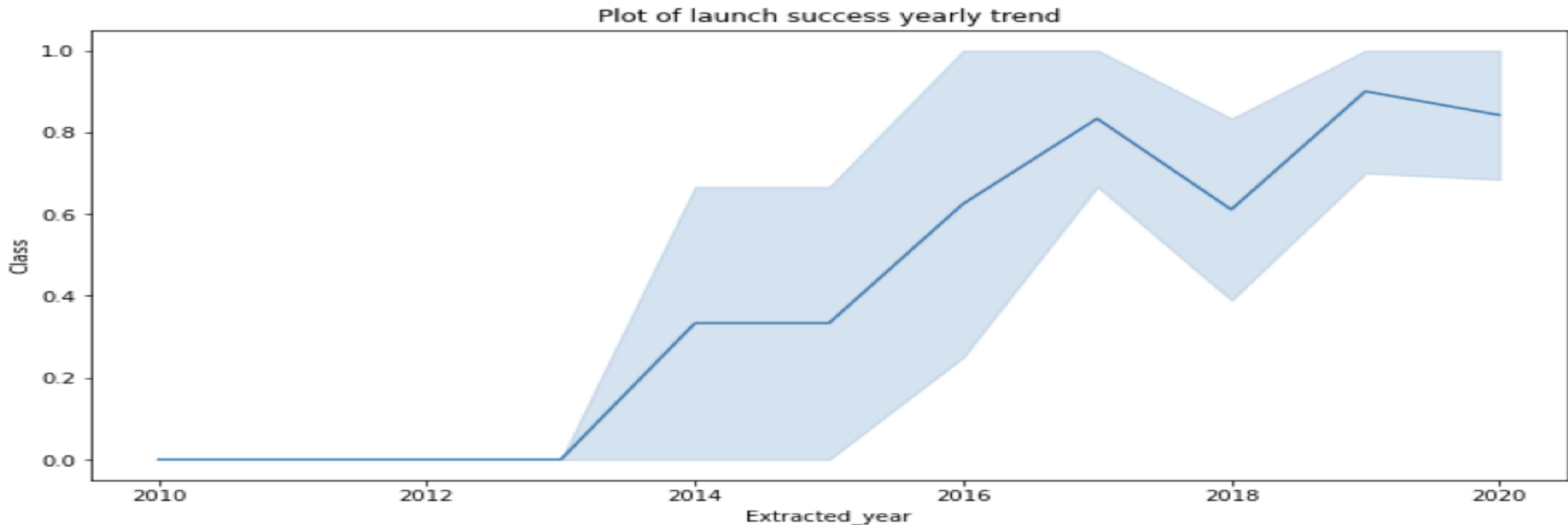
# Success Rate vs. Orbit Type

- From the plot, we can see that ES-L1, GEO, HEO, SSO, VLEO had the most success rate.



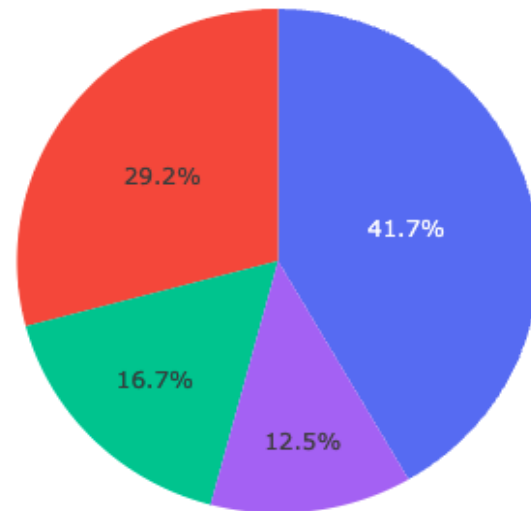
# Launch Success Yearly Trend

- From the plot, we can observe that success rate since 2013 kept on increasing till 2020.



# Pie chart showing the success percentage achieved by each launch site

Success Count for all launch sites



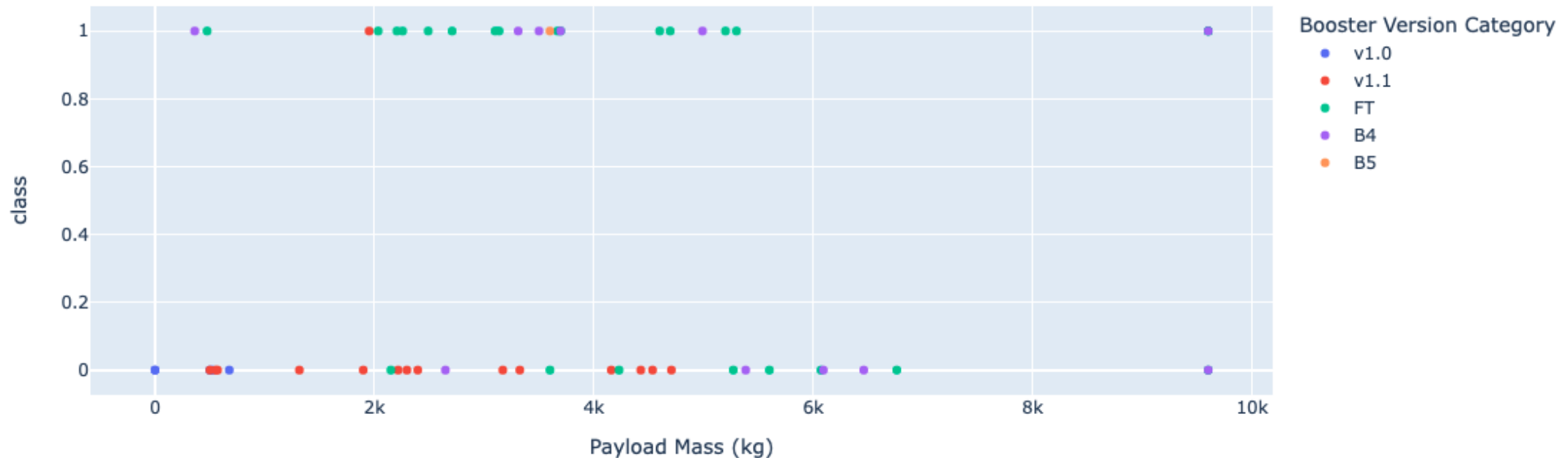
- KSC LC-39A
- CCAFS LC-40
- VAFB SLC-4E
- CCAFS SLC-40

# Scatter plot of Payload mass for all sites

Payload range (Kg):



Success count on Payload mass for all sites



# Classification Accuracy

- The decision tree classifier is the model with the highest classification accuracy. Its score was 0.90

```
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 0.9027777777777778

Best params is : {'criterion': 'entropy', 'max\_depth': 8, 'max\_features': 'auto', 'min\_samples\_leaf': 1, 'min\_samples\_split': 5, 'splitter': 'best'}

# Conclusion

---

- We can conclude that:
- The larger the flight amount at a launch site, the greater the success rate at a launch site.
- Launch success rate started to increase in 2013 till 2020.
- Orbits ES-L1, GEO, HEO, SSO, VLEO had the most success rate.
- KSC LC-39A had the most successful launches of any sites.
- The Decision tree classifier is the best machine learning algorithm for this task.