



IBM Developer  
SKILLS NETWORK

# Winning Space Race with Data Science

Yunus Emre Türker  
05-11-2023



# Outline

---

- Executive Summary
- Introduction
- Methodology
- Results
- Conclusion
- Appendix

# Executive Summary

---

- The following methodologies were used to analyze data:
  - Data Collection using web scraping and SpaceX API;
  - Exploratory Data Analysis (EDA), including data wrangling, data visualization and interactive visual analytics;
  - Machine Learning Prediction.
- Summary of all results
  - It was possible to collect valuable data from public sources;
  - EDA allowed to identify which features are the best to predict success of launchings;
  - Machine Learning Prediction showed the best model to predict which characteristics are important to drive this opportunity by the best way, using all collected data.

# Introduction

---

- The objective is to evaluate the viability of the new company Space Y to compete with Space X.
- Desirable answers:
  - The best way to estimate the total cost for launches, by predicting successful landings of the first stage of rockets;
  - Where is the best place to make launches.



Section 1

# Methodology

# Methodology

---

## Executive Summary

- Data collection methodology:
  - Data from Space X was obtained from 2 sources:
    - Space X API (<https://api.spacexdata.com/v4/rockets/>)
    - WebScraping  
([https://en.wikipedia.org/wiki/List\\_of\\_Falcon/\\_9/\\_and\\_Falcon\\_Heavy\\_launches](https://en.wikipedia.org/wiki/List_of_Falcon/_9/_and_Falcon_Heavy_launches))
- Perform data wrangling
  - Collected data was enriched by creating a landing outcome label based on outcome data after summarizing and analyzing features
- Perform exploratory data analysis (EDA) using visualization and SQL

# Methodology

---

## Executive Summary

- Perform interactive visual analytics using Folium and Plotly Dash
- Perform predictive analysis using classification models
  - Data that was collected until this step were normalized, divided in training and test data sets and evaluated by four different classification models, being the accuracy of each model evaluated using different combinations of parameters.

# Data Collection

---

- Data sets were collected from Space X API (<https://api.spacexdata.com/v4/rockets/>) and from Wikipedia ([https://en.wikipedia.org/wiki/List\\_of\\_Falcon\\_9\\_and\\_Falcon\\_Heavy\\_launches](https://en.wikipedia.org/wiki/List_of_Falcon_9_and_Falcon_Heavy_launches)), using web scraping technics.



# Data Collection - SpaceX API

---

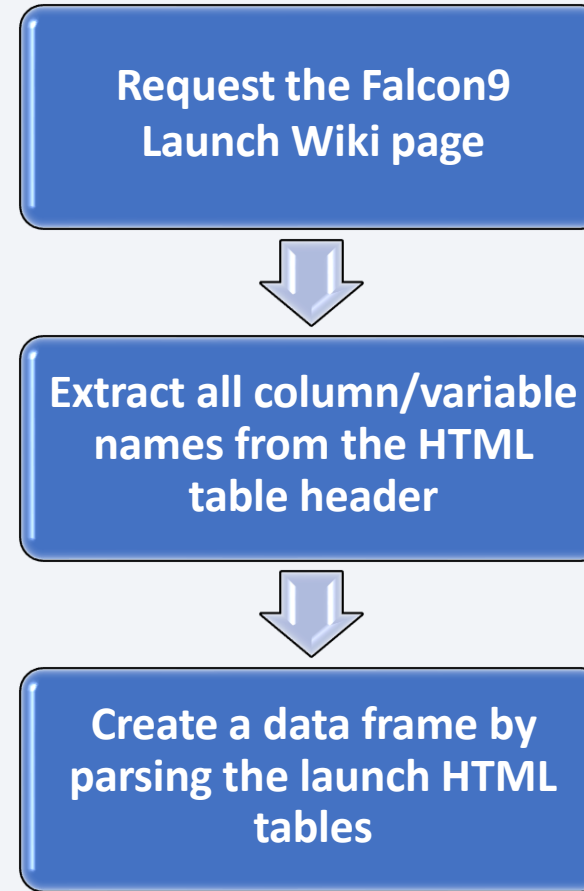
- SpaceX offers a public API from where data can be obtained and then used;
- This API was used according to the flowchart beside and then data is persisted.
- <https://github.com/yunussturker/testrepo/blob/main/Data%20Collection%20with%20Web%20Scraping.ipynb>



# Data Collection - Scraping

---

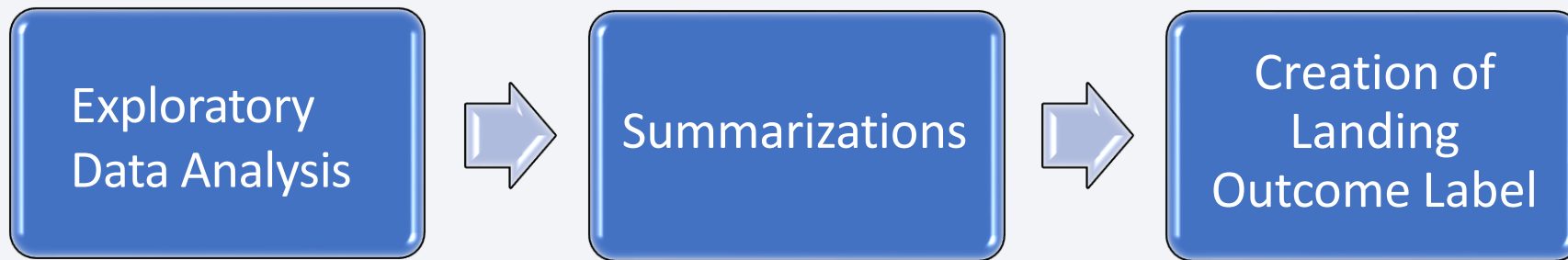
- Data from SpaceX launches can also be obtained from Wikipedia;
- Data are downloaded from Wikipedia according to the flowchart and then persisted.



# Data Wrangling

---

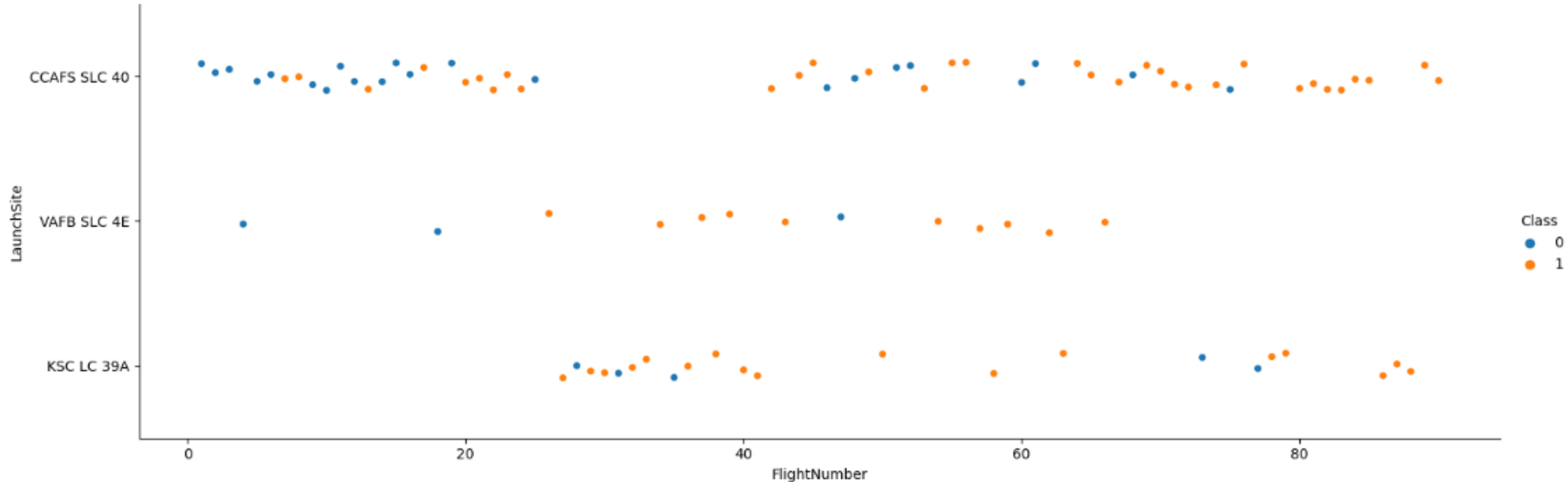
- Initially some Exploratory Data Analysis was performed on the dataset.
- Then the summaries launches per site, occurrences of each orbit and occurrences of mission outcome per orbit type were calculated.
- Finally, the landing outcome label was created from Outcome column.



- <https://github.com/yunussturker/testrepo/blob/main/Data%20Wrangling.ipynb>

# EDA with Data Visualization

```
### TASK 1: Visualize the relationship between Flight Number and Launch Site  
sns.catplot(y="LaunchSite", x="FlightNumber", hue="Class", data=df, aspect = 3)  
plt.show()
```



- <https://github.com/yunussturker/testrepo/blob/main/DataVisualization.ipynb>

# EDA with SQL

---

- The following SQL queries were performed:
  - Names of the unique launch sites in the space mission;
  - Top 5 launch sites whose name begin with the string 'CCA';
  - Total payload mass carried by boosters launched by NASA (CRS);
  - Average payload mass carried by booster version F9 v1.1;
  - Date when the first successful landing outcome in ground pad was achieved;
  - Names of the boosters which have success in drone ship and have payload mass between 4000 and 6000 kg;
  - Total number of successful and failure mission outcomes;
  - Names of the booster versions which have carried the maximum payload mass;
  - Failed landing outcomes in drone ship, their booster versions, and launch site names for in year 2015; and
  - Rank of the count of landing outcomes (such as Failure (drone ship) or Success (ground pad)) between the date 2010-06-04 and 2017-03-20.
- <https://github.com/yunussturker/testrepo/blob/main/ExploratoryDataAnalysis.ipynb>

# Build an Interactive Map with Folium

---

- Markers, circles, lines and marker clusters were used with Folium Maps
  - Markers indicate points like launch sites;
  - Circles indicate highlighted areas around specific coordinates, like NASA Johnson Space Center;
  - Marker clusters indicates groups of events in each coordinate, like launches in a launch site; and
  - Lines are used to indicate distances between two coordinates.
- <https://github.com/yunussturker/testrepo/blob/main/FoliumVisualization.ipynb>



# Build a Dashboard with Plotly Dash

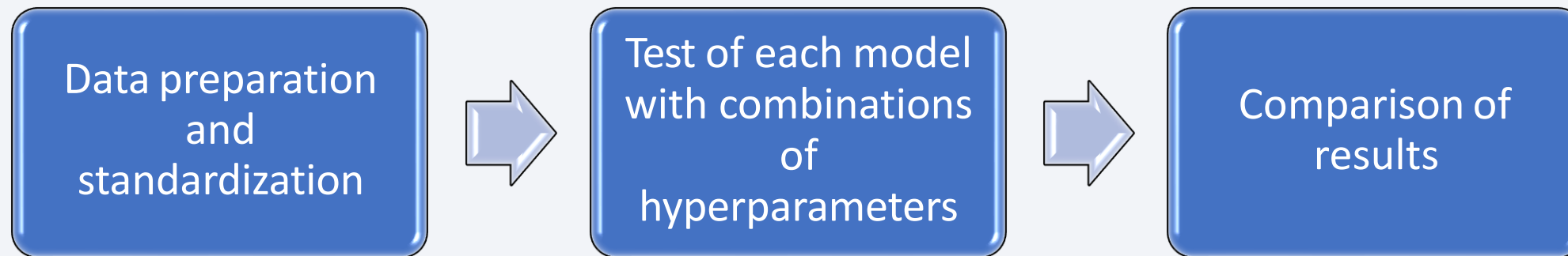
---

- The following graphs and plots were used to visualize data
  - Percentage of launches by site
  - Payload range
- This combination allowed to quickly analyze the relation between payloads and launch sites, helping to identify where is best place to launch according to payloads.
- <https://github.com/yunussturker/testrepo/blob/main/DashVisualization.py>

# Predictive Analysis (Classification)

---

- Four classification models were compared: logistic regression, support vector machine, decision tree and k nearest neighbors.



- <https://github.com/yunussturker/testrepo/blob/main/MachineLearning.ipynb>

# Results

---

- Exploratory data analysis results:
  - Space X uses 4 different launch sites;
  - The first launches were done to Space X itself and NASA;
  - The average payload of F9 v1.1 booster is 2,928 kg;
  - The first success landing outcome happened in 2015 five year after the first launch;
  - Many Falcon 9 booster versions were successful at landing in drone ships having payload above the average;
  - Almost 100% of mission outcomes were successful;
  - Two booster versions failed at landing in drone ships in 2015: F9 v1.1 B1012 and F9 v1.1 B1015;
  - The number of landing outcomes became as better as years passed.

# Results

- Using interactive analytics was possible to identify that launch sites use to be in safety places, near sea, for example and have a good logistic infrastructure around.
- Most launches happens at east cost launch sites.



# Results

---

- Predictive Analysis showed that all classification algorithms have same accuracy for this working.

## TASK 12

Find the method performs best:

```
print("Logistic Regression Accuracy: ", logreg_cv.score(X_test, y_test))
print("Decision Tree Accuracy:      ", tree_cv.score(X_test, y_test))
print("SVM Accuracy:                ", svm_cv.score(X_test, y_test))
print("KNN Accuracy:                ", knn_cv.score(X_test, y_test))
```

```
Logistic Regression Accuracy: 0.8333333333333334
Decision Tree Accuracy:      0.8333333333333334
SVM Accuracy:                0.8333333333333334
KNN Accuracy:                0.8333333333333334
```



The background of the slide is an abstract composition. It features a solid blue area on the left side, which transitions into a dynamic pattern of diagonal streaks in shades of blue, red, and cyan on the right. These streaks vary in thickness and intensity, creating a sense of motion and depth. A faint, light-blue grid pattern is visible across the entire background, adding a technical or digital feel to the design.

Section 2

# Insights drawn from EDA



# Flight Number vs. Launch Site



- According to the plot above, it's possible to verify that the best launch site nowadays is CCAFS SLC 40, where most of recent launches were successful;
- In second place VAFB SLC 4E and third place KSC LC 39A;
- It's also possible to see that the general success rate improved over time.

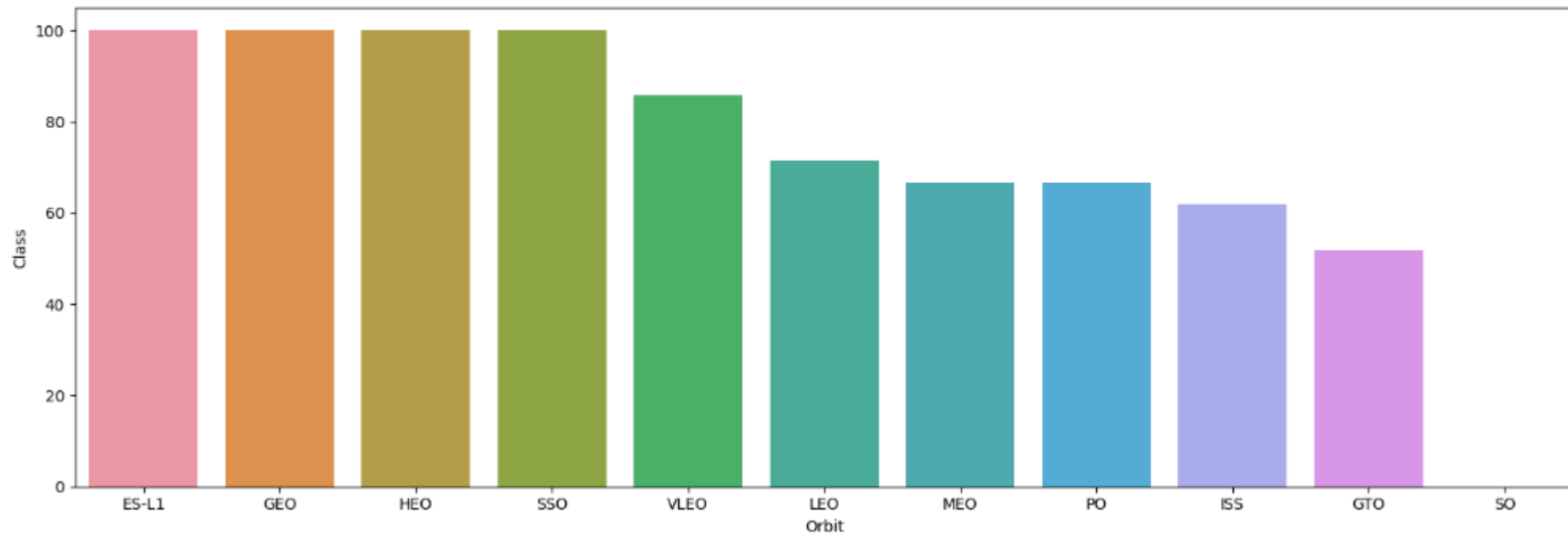
# Payload vs. Launch Site



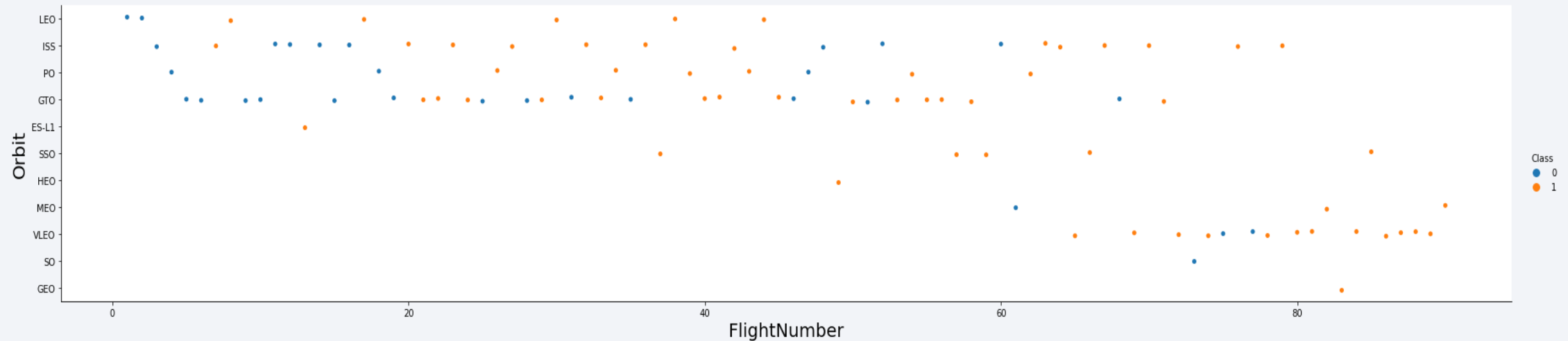
- Payloads over 9,000kg (about the weight of a school bus) have excellent success rate;
- Payloads over 12,000kg seems to be possible only on CCAFS SLC 40 and KSC LC 39A launch sites.

# Success Rate vs. Orbit Type

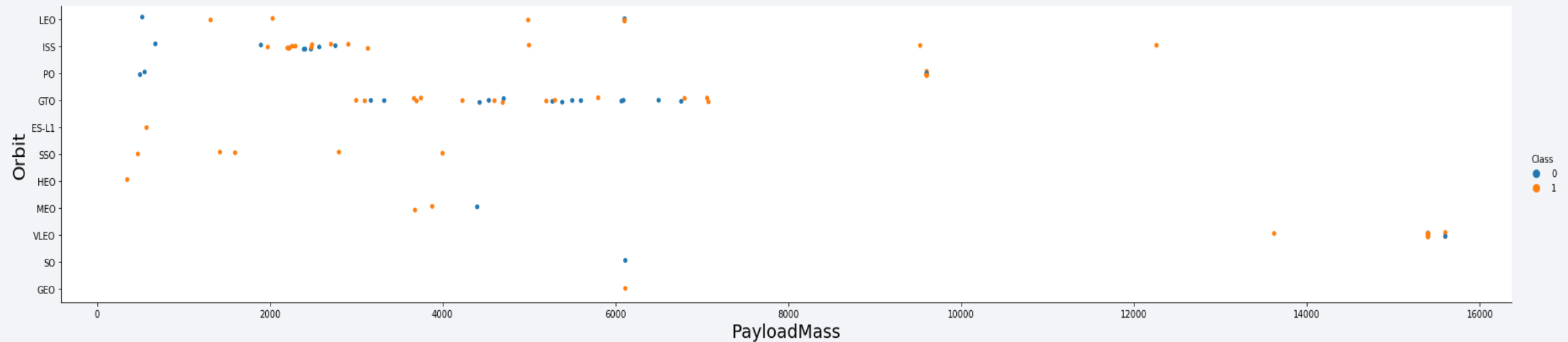
```
### TASK 3: Visualize the relationship between success rate of each orbit type
group = df.groupby("Orbit")["Class"].mean().reset_index().sort_values(by="Class", ascending=False)
group["Class"] = group["Class"] * 100
sns.barplot(data=group, x="Orbit", y="Class")
plt.show()
```



# Flight Number vs. Orbit Type



# Payload vs. Orbit Type

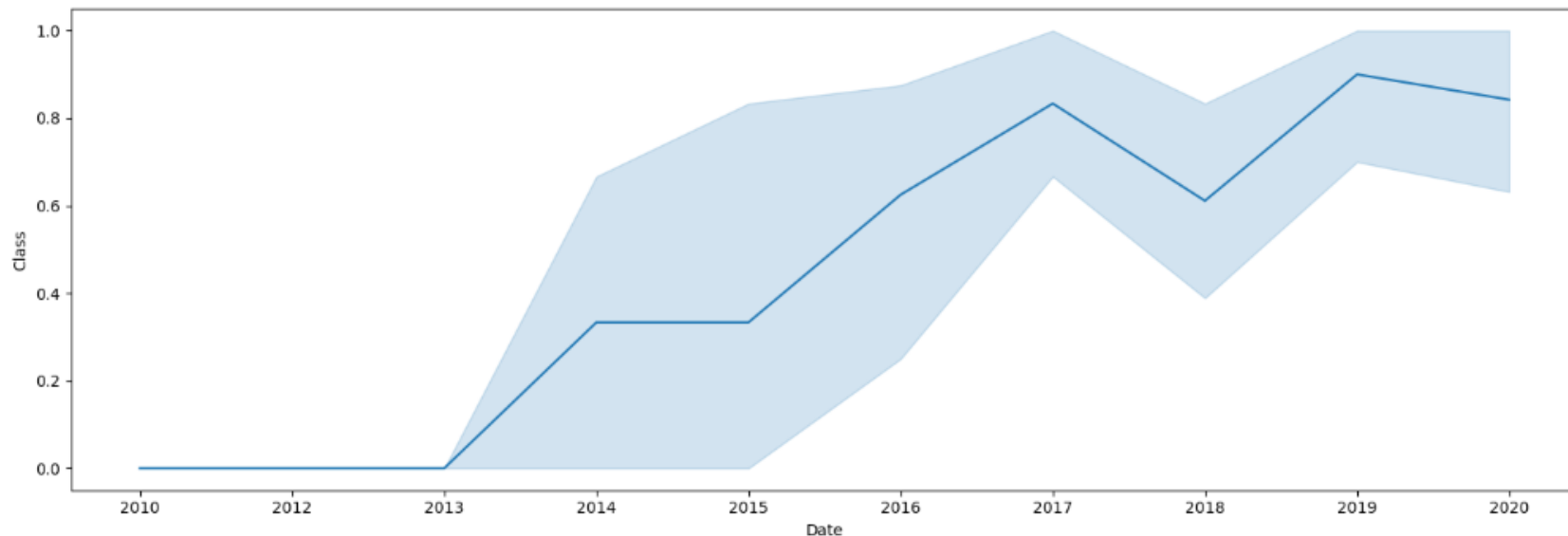


- Apparently, there is no relation between payload and success rate to orbit GTO;
- ISS orbit has the widest range of payload and a good rate of success;
- There are few launches to the orbits SO and GEO.

# Launch Success Yearly Trend

---

- Success rate started increasing in 2013 and kept until 2020;
- It seems that the first three years were a period of adjusts and improvement of technology.





# All Launch Site Names

---

- According to data, there are four launch sites:

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

- They are obtained by selecting unique occurrences of “launch\_site” values from the dataset.

# Launch Site Names Begin with 'CCA'

---

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

Launch_Site
CCAFS LC-40
CCAFS LC-40
CCAFS LC-40
CCAFS LC-40
CCAFS LC-40

- Here we can see five samples of Cape Canaveral launches.

# Total Payload Mass

---

- Total payload carried by boosters from NASA:

SUM(PAYLOAD_MASS_KG_)
107010

- Total payload calculated above, by summing all payloads whose codes contain 'CRS', which corresponds to NASA.

# Average Payload Mass by F9 v1.1

---

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

```
AVG(PAYLOAD_MASS__KG_)
2534.6666666666665
```

- Filtering data by the booster version above and calculating the average payload mass we obtained the value of 2534.66 kg.

# First Successful Ground Landing Date

---

- First successful landing outcome on ground pad:

Date	Landing _Outcome
22-12-2015	Success (ground pad)

- By filtering data by successful landing outcome on ground pad and getting the minimum value for date it's possible to identify the first occurrence, that happened on 22-12-2015.

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

---

- Boosters which have successfully landed on drone ship and had payload mass greater than 4000 but less than 6000

Booster_Version	PAYLOAD_MASS_KG_	Mission_Outcome	Landing_Outcome
F9 FT B1022	4696	Success	Success (drone ship)
F9 FT B1026	4600	Success	Success (drone ship)
F9 FT B1021.2	5300	Success	Success (drone ship)
F9 FT B1031.2	5200	Success	Success (drone ship)

- Selecting distinct booster versions according to the filters above, these 4 are the result.



# Total Number of Successful and Failure Mission Outcomes

---

- Number of successful and failure mission outcomes:

Mission_Outcome	COUNT(Mission_Outcome)
Failure (in flight)	1
Success	98
Success	1
Success (payload status unclear)	1

- Grouping mission outcomes and counting records for each group led us to the summary above.

# Boosters Carried Maximum Payload

---

- Boosters which have carried the maximum payload mass

Booster_Version	PAYLOAD_MASS_KG_
F9 B5 B1048.4	15600
F9 B5 B1049.4	15600
F9 B5 B1051.3	15600
F9 B5 B1056.4	15600
F9 B5 B1048.5	15600
F9 B5 B1051.4	15600
F9 B5 B1049.5	15600
F9 B5 B1060.2	15600
F9 B5 B1058.3	15600
F9 B5 B1051.6	15600
F9 B5 B1060.3	15600
F9 B5 B1049.7	15600

- These are the boosters which have carried the maximum payload mass registered in the dataset.

# 2015 Launch Records

---

- Failed landing outcomes in drone ship, their booster versions, and launch site names for in year 2015

MONTH	Booster_Version	Launch_Site	Landing_Outcome
01	F9 v1.1 B1012	CCAFS LC-40	Failure (drone ship)
04	F9 v1.1 B1015	CCAFS LC-40	Failure (drone ship)

- The list above has the only two occurrences.

# Rank Landing Outcomes Between 2020-06-04 and 2017-03-20

---

- Ranking of all landing outcomes between the date 2020-06-04 and 2017-03-20:

Date	Landing _Outcome
06-12-2020	Success
16-11-2020	Success
05-11-2020	Success
18-10-2020	Success
06-10-2020	Success
18-08-2020	Success
07-08-2020	Success
13-06-2020	Success
04-06-2020	Success
07-03-2020	Success
07-01-2020	Success
17-12-2019	Success
05-12-2019	Success
11-11-2019	Success
12-06-2019	Success
11-01-2019	Success
15-11-2018	Success
08-10-2018	Success
10-09-2018	Success
07-08-2018	Success

Section 4

# Launch Sites Proximities Analysis



# All launch sites

---



- Launch sites are near sea, probably by safety, but not too far from roads and railroads.

# Launch Outcomes by Site

- Example of KSC LC-39A launch site launch outcomes



- Green markers indicate successful and red ones indicate failure.



# Logistics and Safety

---



- Launch site KSCLC-39A has good logistics aspects, being near railroad and road and relatively far from inhabited areas.

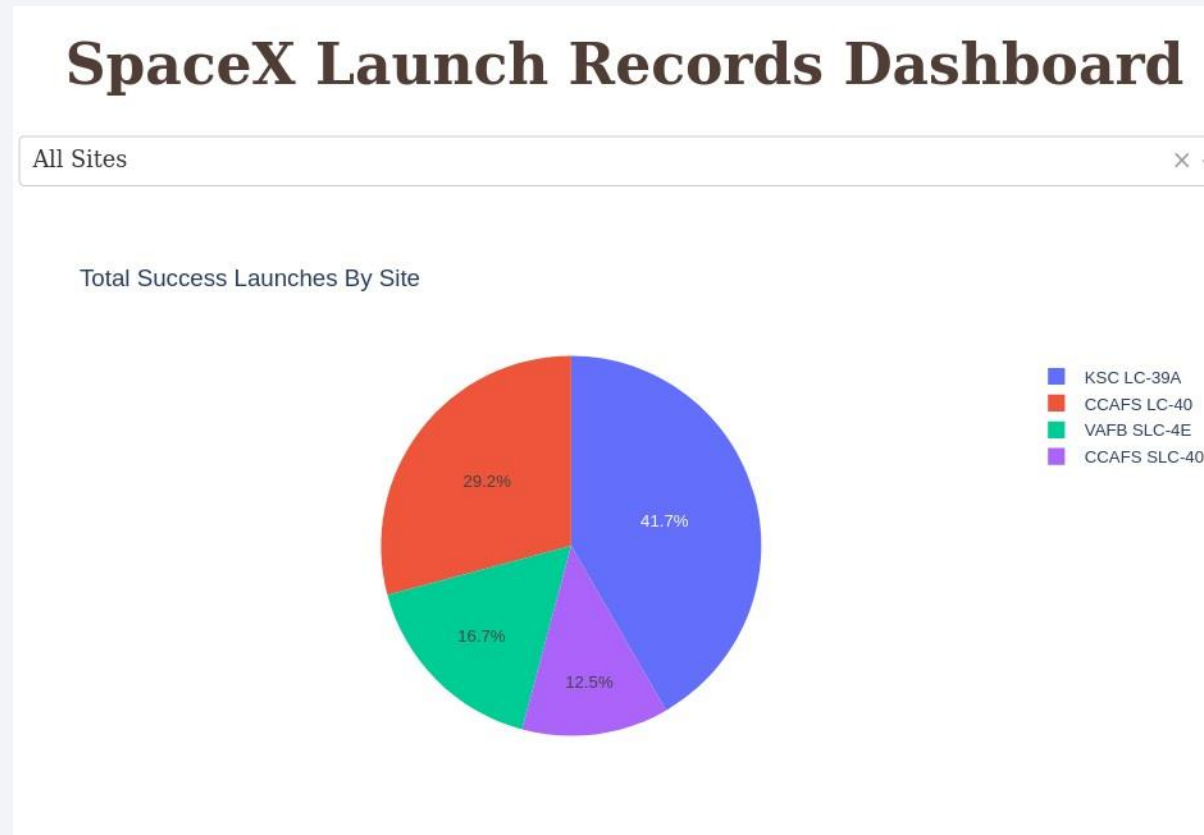




Section 5

# Build a Dashboard with Plotly Dash

# Successful Launches by Site



- The place from where launches are done seems to be a very important factor of success of missions.

# Launch Success Ratio for KSC LC-39A

---



- 76.9% of launches are successful in this site.

# Payload vs. Launch Outcome



- Payloads under 6,000kg and FT boosters are the most successful combination.

# Payload vs. Launch Outcome



- There's not enough data to estimate risk of launches over 7,000kg

Section 6

# Predictive Analysis (Classification)



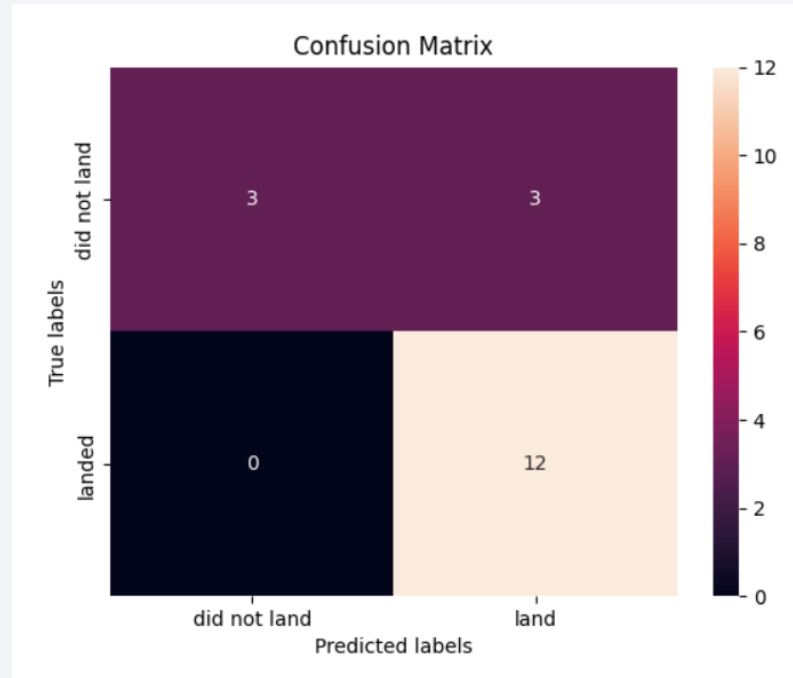
# Classification Accuracy

---

- Four classification models were tested, and their accuracies are plotted beside;
- The model with the highest classification accuracy all models have same accuracy.

```
Logistic Regression Accuracy: 0.8333333333333334  
Decision Tree Accuracy:      0.8333333333333334  
SVM Accuracy:                0.8333333333333334  
KNN Accuracy:                0.8333333333333334
```

# Confusion Matrix of Decision Tree Classifier



- Confusion matrix of Decision Tree Classifier proves its accuracy by showing the big numbers of true positive and true negative compared to the false ones.



# Conclusions

---

- Different data sources were analyzed, refining conclusions along the process;
- The best launch site is KSC LC-39A;
- Launches above 7,000kg are less risky;
- Although most of mission outcomes are successful, successful landing outcomes seem to improve over time, according the evolution of processes and rockets;
- Decision Tree Classifier can be used to predict successful landings and increase profits.

# Appendix

---

- As an improvement for model tests, it's important to set a value to `np.random.seed` **variable**;
- Folium didn't show maps on Github, so I took screenshots.

Thank you!

