# Outline

- Executive Summary

- Introduction

- Methodology

- Results

- Conclusion

- Appendix

# Executive Summary

- This analysis collected data from SpaceX API and Wikipedia and transformed into a dataframe. it was later visualized, queried, and classified

  - It follows normal EDA process, used visualization, SQL, and machine learning technology to proceed.

- Summary of all results

  - From the data analysis, we can find that some factors such as payload mass, orbit, and time are all influential to the success rate

  - Through folium map, we can find some common attributes from the launch sites

  - When we are trying to predict the probability of success, all the methods are doing well

# Introduction

- Since 2002, SpaceX has become an undeniable force in the global aerospace industry with its innovative spirit and relentless pursuit. However, rocket launches, like other scientific explorations, are full of unknowns and uncertainties.

- The purpose of this study is to review the rocket launch records of Falcon 9 and analyze the key factors that affect the launch results in a data science manner. Hopefully, this analysis can start with data and provide some insights for the subsequent launch of SpaceX.

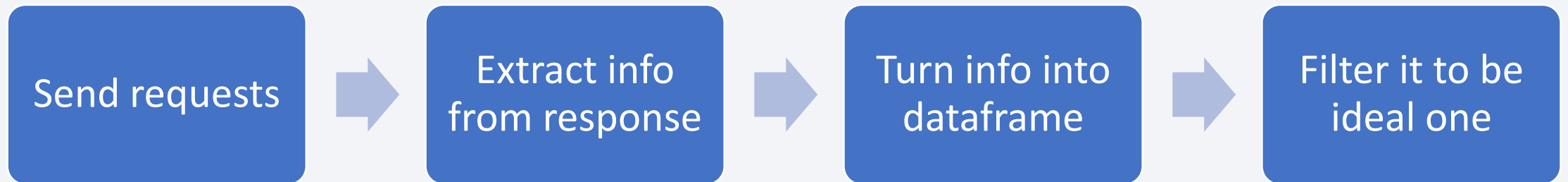Section 1

# Methodology

# Methodology

## Executive Summary

- Data collection methodology:

  - The information is sourced from the SpaceX API and Wikipedia. They were discovered from web pages and parsed from responses in Json and Html formats

- Perform data wrangling

  - Here all information was normalized, checked, and re-organized into a dataframe

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

- Perform interactive visual analytics using Folium and Plotly Dash

- Perform predictive analysis using classification models

  - The dataset was separated into train and test set, they were further fitted by different models.

# Data Collection

- The datasets were collected from SpaceX API and Wikipedia.

- The collection started with requests and getting response in different format, and the response were all transformed into dataframe which is also filtered later.

```
Send requests  →  Extract info from response  →  Turn info into dataframe  →  Filter it to be ideal one
```

# Data Collection – SpaceX API

In this part:

- Requests of rocket launch data had been sent to SpaceX API with URL

- the Json result had been normalized and turned it into a dataframe

- The original dataframe was then filtered with a defined dictionary and further filtered to be a falcon 9 only dataframe

- For detail of the workflow, you may check the notebook here.

# Data Collection - Scraping

In this part:

- Requests of rocket launch data had been sent to Wiki

- Target was parsed with beautiful soup

- A dictionary had been defined to reorganize the information and it was later turned into a dataframe

- For detail of the workflow, you may check the notebook here.

# Data Wrangling

During the process to get the ideal dataframe:

- All null values had been replaced by average number to ensure no abnormalities

- Types of data were checked

- Occurrence of orbits and outcomes were also checked

- The "class" column was created from outcomes, to comprehensively separate good and bad outcomes

- For detail of the workflow, you may check the notebook here.

Replace the nulls

Check data

Re-organize the frame

# EDA with Data Visualization

- Scatter was used here to estimate the relationship between different factors

    - FlightNumber vs. Launch Site: relationship between Flight Number and Launch Site

    - Payload Mass vs. Launch Site: relationship between Payload Mass and Launch Site

    - Flight Number vs. Orbit Type: relationship between Flight Number and Orbit type

    - Payload vs. Orbit Type: relationship between Payload and Orbit type

- Success Rate vs. Orbit Type: a bar chart to compare the success rates of each orbit

- A Line was also plotted to show the Launch Success Yearly Trend:

- Further detail, please check the notebook here.

# EDA with SQL

- SQL queries performed to this data set

  - select distinct Launch_Site from SPACEXTABLE

  - select * from SPACEXTABLE WHERE Launch_Site LIKE 'CCA%' LIMIT 5

  - select sum(PAYLOAD_MASS__KG_) from SPACEXTABLE WHERE Customer='NASA (CRS)'

  - select avg(PAYLOAD_MASS__KG_) from SPACEXTABLE WHERE Booster_Version='F9 v1.1'

  - select min(Date) from SPACEXTABLE where Landing_Outcome='Success (ground pad)'

  - select * from SPACEXTABLE where Landing_Outcome='Success (drone ship)' AND PAYLOAD_MASS__KG_>4000 AND PAYLOAD_MASS__KG_<6000

  - select Mission_Outcome, count(*) from SPACEXTABLE group by Mission_Outcome

  - select Booster_Version from SPACEXTABLE where PAYLOAD_MASS__KG_=(select max(PAYLOAD_MASS__KG_) from SPACEXTABLE)

  - select substr(Date, 6,2) as Month, Landing_Outcome, Booster_Version, Launch_Site from SPACEXTABLE where Landing_Outcome='Failure (drone ship)' and substr(Date,0,5)='2015'

  - select Landing_Outcome, COUNT(*) AS Outcome_No from SPACEXTABLE where Date between '2010-06-04' and '2017-03-20' group by Landing_Outcome order by Outcome_No DESC

- Result of queries, please check [here](here)                12

# Build an Interactive Map with Folium

- Several objects have been created and added to a folium map, including:

  - Circles and markers for all launch sites on the map

  - Markers for the success/failed launches for each site on the map

  - Line for the closest coast line

- This is for you to:

  - Notice the location of all sites, and it will enable you to select the site and read its description

  - Find a cluster of outcomes from different sites, which you may click to see every outcomes happened

  - Understand the selection of those locations for safety concern of a bad outcome

- Further detail, please check [here](here)

# Build a Dashboard with Plotly Dash

- Here are 5 pies and 5 scatters in this dashboard:

    - If you select all sites(default), it will be the share of launches from 4 sites, and the scatter down below will contain all samples

    - If you select a launch site, it will return you a the success share of this site, and the scatter down below will contain the samples from this site

    - You can also use the slide bar to select the range of payload


- For the code, please check [here](#)

# Predictive Analysis (Classification)

- In order to predict the probability of a good or bad outcome, column "class" was set to Y, other columns were set to X

- 20% of the data set was set as test group

- 4 models have been applied to this dataset

  - Logic regression

  - SVM

  - Decision tree

  - KNN

- Score function and confusion matrix were applied to evaluate the models

- Eventually, all models returned the same result in test

- For the full process, please check [here](#)

# Results

- Exploratory data analysis results

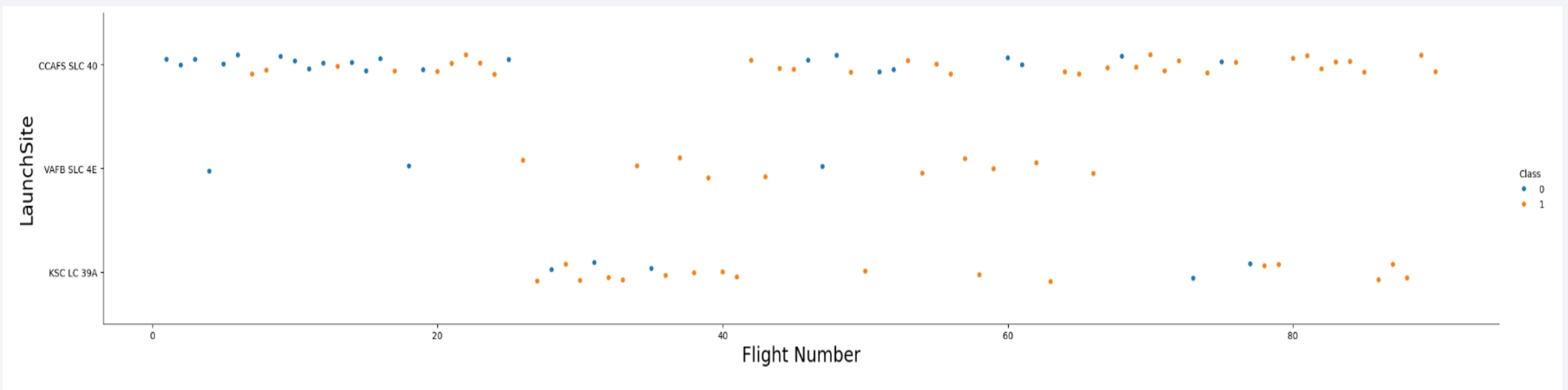- Interactive analytics demo in screenshots

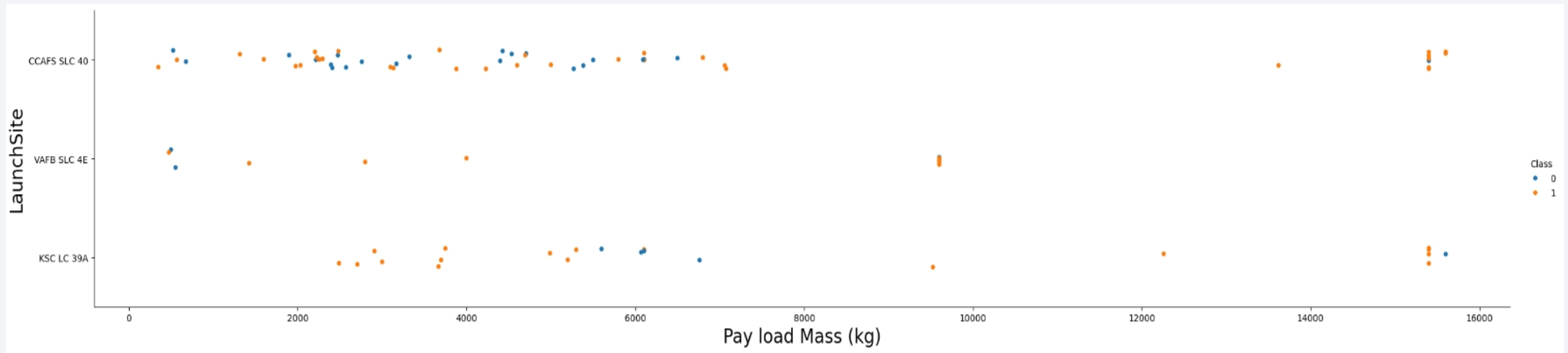- Predictive analysis results

Section 2

# Insights drawn from EDA

# Flight Number vs. Launch Site

- VAFB SLC 4E has the best success rate

- However, it may because of its small sample

- As a result, we cannot say the outcome is relevant with launch site yet

# Payload vs. Launch Site

- Replacing the flight number with payload can significantly bring some change

- Obviously, VAFB SLC 4E has several outcomes from the same payload which is in a high success rate area

- Compare with the launch site, the payload seems to be a bigger factor

- However, though the sample size is still small, KSC LC 39A and VAFB 4E showed better potential than CCAFS SLC40 when payload is between 1500 and 5500

# Success Rate vs. Orbit Type

- Similar process can been repeated in Orbit types

- A bar chart can pick the Orbits with better success rate immediately

# Flight Number vs. Orbit Type

- The flight number shows that several Orbits have significant less samples

- Most 100% success Orbits has only one outcome

- SSO is the only 100% success Orbit with more than one outcome, but still a small subset

# Payload vs. Orbit Type

- All SSO samples are from the payload area 0-6000, where you may still remember a lot of failure from CCAFS site

- LEO is not too far from SSO with only twice bad outcome in the same payload range

- Still, payload > 6000 range, not many bad outcomes

# Launch Success Yearly Trend

- In flight number scatter, a noticeable fact is: bigger number works better

- The flight number is relevant to the time and experience, and this guess is verified by the trend line right side

# All Launch Site Names

- Through SQL query, here you may find the unique members of Launch_Site family. In another word, we got 4 Launch Sites in total and here they are.

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

# Launch Site Names Begin with 'CCA'

- There are 2 launch sites starts with 'CCA'

- They are filtered out

| Date | Time (UTC) | Booster_Version | Launch_Site | Payload | PAYLOAD_MASS__KG_ | Orbit | Customer | Mission_Outcome | Landing_Outcome |
|---|---|---|---|---|---|---|---|---|---|
| 2010-06-04 | 18:45:00 | F9 v1.0 B0003 | CCAFS LC-40 | Dragon Spacecraft Qualification Unit | 0 | LEO | SpaceX | Success | Failure (parachute) |
| 2010-12-08 | 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 | 7:44:00 | F9 v1.0 B0005 | CCAFS LC-40 | Dragon demo flight C2 | 525 | LEO (ISS) | NASA (COTS) | Success | No attempt |
| 2012-10-08 | 0:35:00 | F9 v1.0 B0006 | CCAFS LC-40 | SpaceX CRS-1 | 500 | LEO (ISS) | NASA (CRS) | Success | No attempt |
| 2013-03-01 | 15:10:00 | F9 v1.0 B0007 | CCAFS LC-40 | SpaceX CRS-2 | 677 | LEO (ISS) | NASA (CRS) | Success | No attempt |

# Total Payload Mass

- the total payload carried by boosters from NASA is 45596kg

**sum(PAYLOAD_MASS__KG_)**

45596

# Average Payload Mass by F9 v1.1

- The average payload mass carried by booster version F9 v1.1 is 2928.4kg

**avg(PAYLOAD_MASS__KG_)**

2928.4

# First Successful Ground Landing Date

- The dates of the first successful landing outcome on ground pad is 2015/12/22

**min(Date)**

2015-12-22

# Successful Drone Ship Landing with Payload between 4000 and 6000

- 4 records have been found landed on a drone ship with payload 4000-6000 kg

- 2 from CCAFS LC-40 with payload around 4650, 2 from KSC LC-39A with payload around 5250

- All for GTO

| Date | Time (UTC) | Booster_ Version | Launch_Site | Payload | PAYLOAD_ MASS__KG_ | Orbit | Customer | Mission_ Outcome | Landing_ Outcome |
|---|---|---|---|---|---|---|---|---|---|
| 2016-05-06 | 5:21:00 | F9 FT B1022 | CCAFS LC-40 | JCSAT-14 | 4696 | GTO | SKY Perfect JSAT Group | Success | Success (drone ship) |
| 2016-08-14 | 5:26:00 | F9 FT B1026 | CCAFS LC-40 | JCSAT-16 | 4600 | GTO | SKY Perfect JSAT Group | Success | Success (drone ship) |
| 2017-03-30 | 22:27:00 | F9 FT B1021.2 | KSC LC-39A | SES-10 | 5300 | GTO | SES | Success | Success (drone ship) |
| 2017-10-11 | 22:53:00 | F9 FT B1031.2 | KSC LC-39A | SES-11 / EchoStar 105 | 5200 | GTO | SES EchoStar | Success | Success (drone ship) |

# Total Number of Successful and Failure Mission Outcomes

- Almost all the boosters recorded as success in mission outcome category

- However mission success doesn't mean successfully landed

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

# Boosters Carried Maximum Payload

- Here we got the booster carried the maximum payload mass

**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

- Here is the ones failed to land in drone ship in 2015

- Both of them are from 1H and Booster Version is reasonably close

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

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

- From 2010-06-04 to 2017-03-20, the most landing outcome is "no attempt" with 10 times

- For drone ship, its landing outcome is evenly distributed as 5

- Boosters successfully landed on the ground pad for 3 times

| Landing_Outcome | Outcome_No |
|---|---|
| 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 |

Section 3

# Launch Sites Proximities Analysis

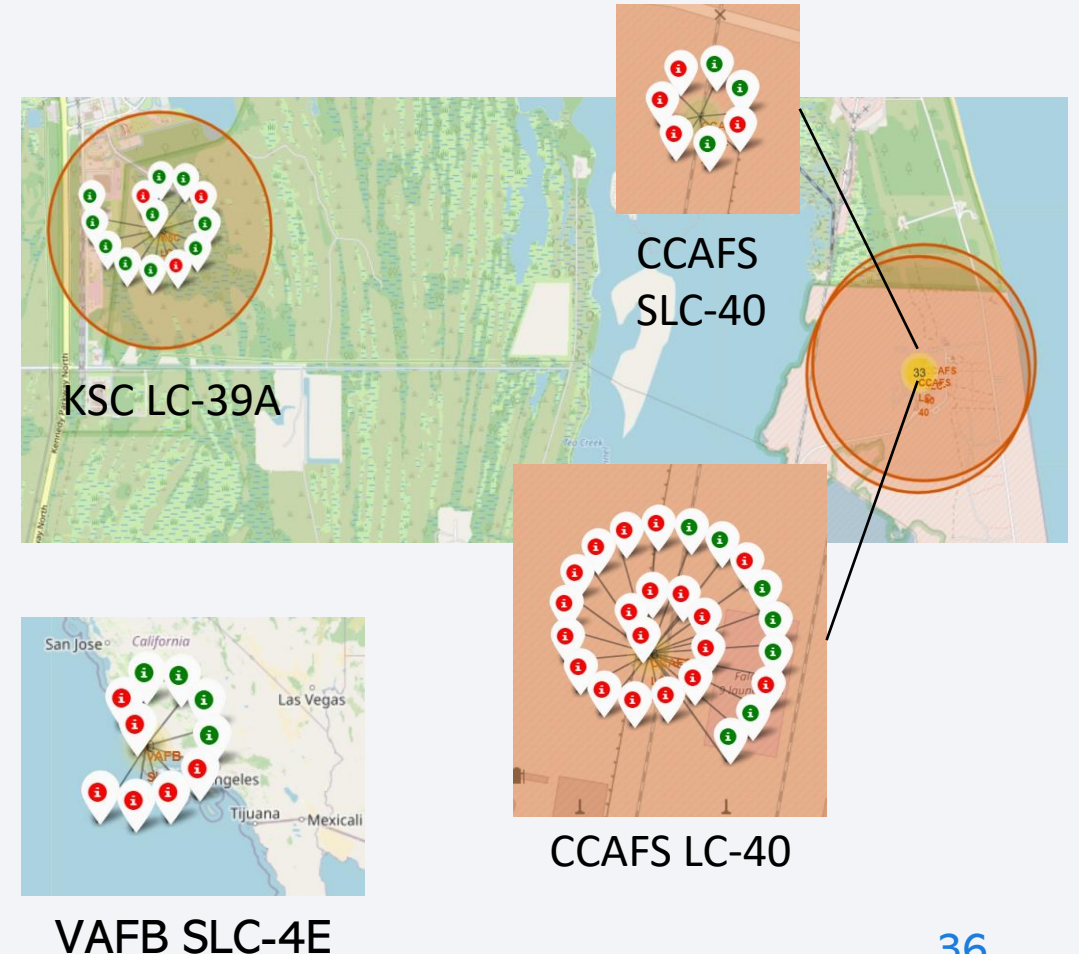# Location of Launch sites

- Through this map, we can easily find some fact:

  - Only 1 site from west coast

  - 3 sites located very closely in Florida
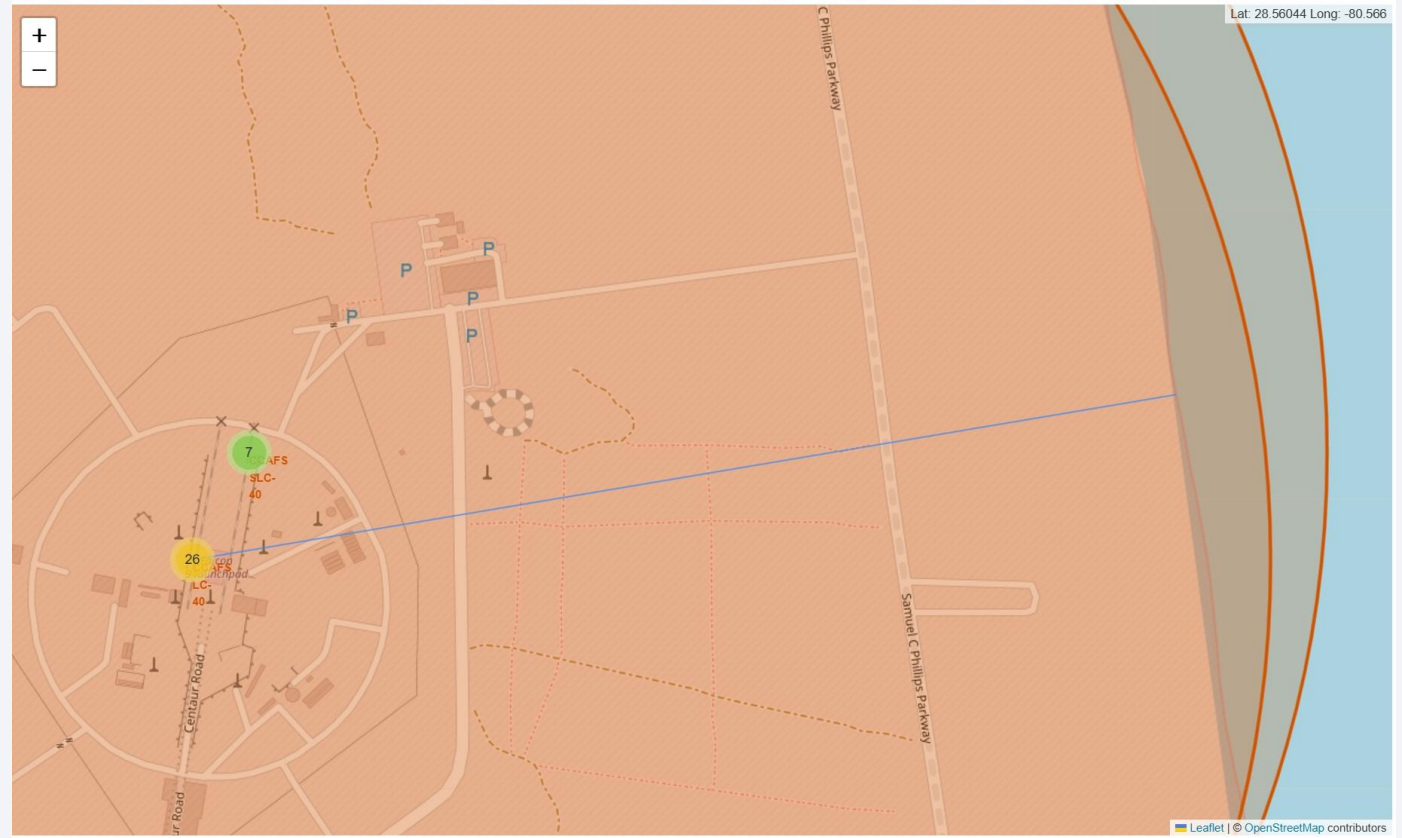
  - All the sites are beside coasts

# Outcomes from each Launch Sites

- KSC LC-39A has the best record

- All other sites have bad outcomes more than good outcomes

- CCAFS SLC-40 located at almost the same place of CCAFS LC-40 but recorded 3 good outcomes and 4 bad outcomes, which is closer to 50%



CCAFS SLC-40

KSC LC-39A

CCAFS LC-40

VAFB SLC-4E

# CCAFS sites are very close to the coast

- Different from KSC LC-39A, both CCAFS sites are located less than 1km from the coast line.

- As previously mentioned, CCAFS sites has the most launch record, also the failed ones.

- Such location offered better chance to turn bad outcome into controlled or uncontrolled (ocean)
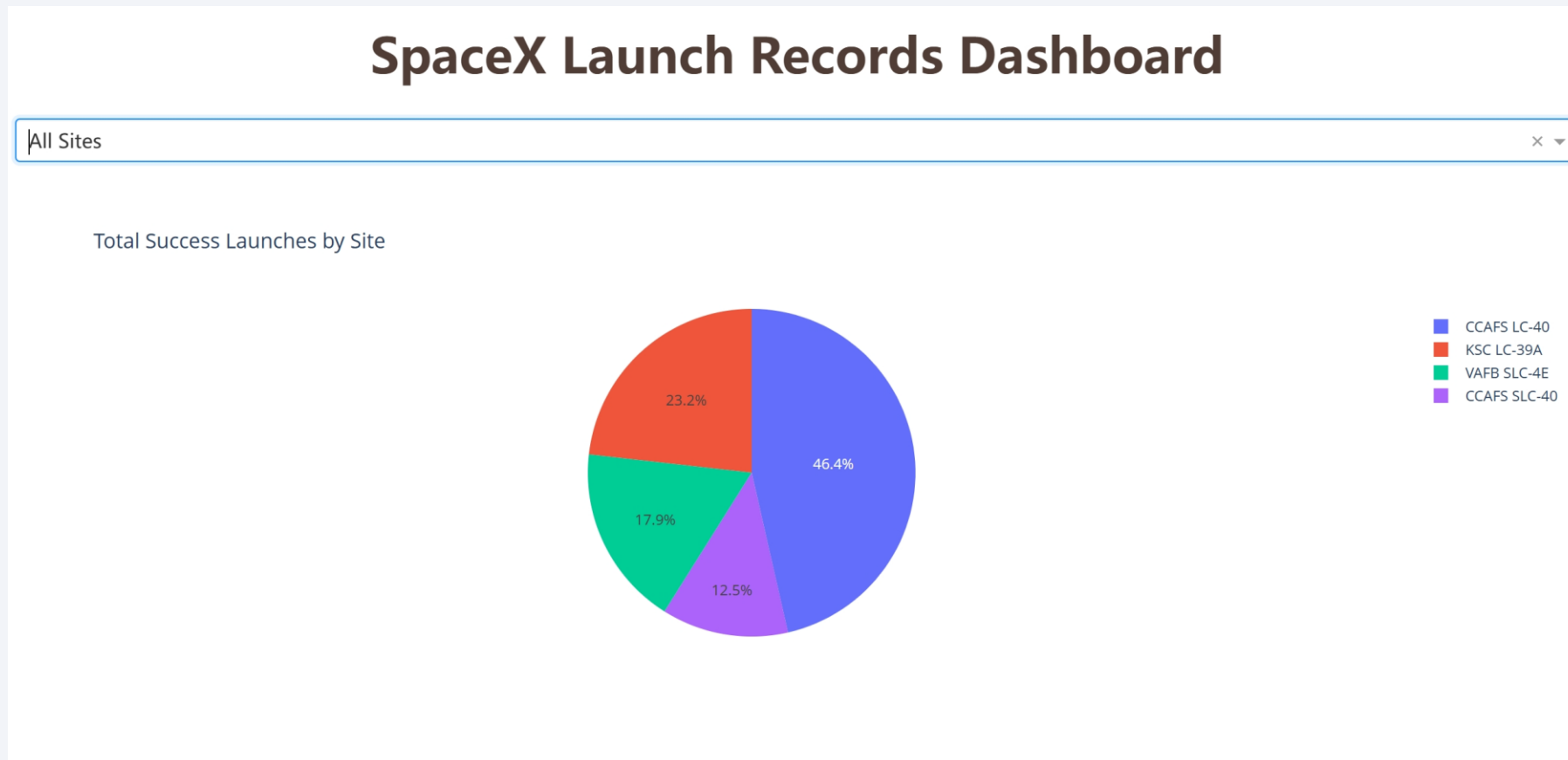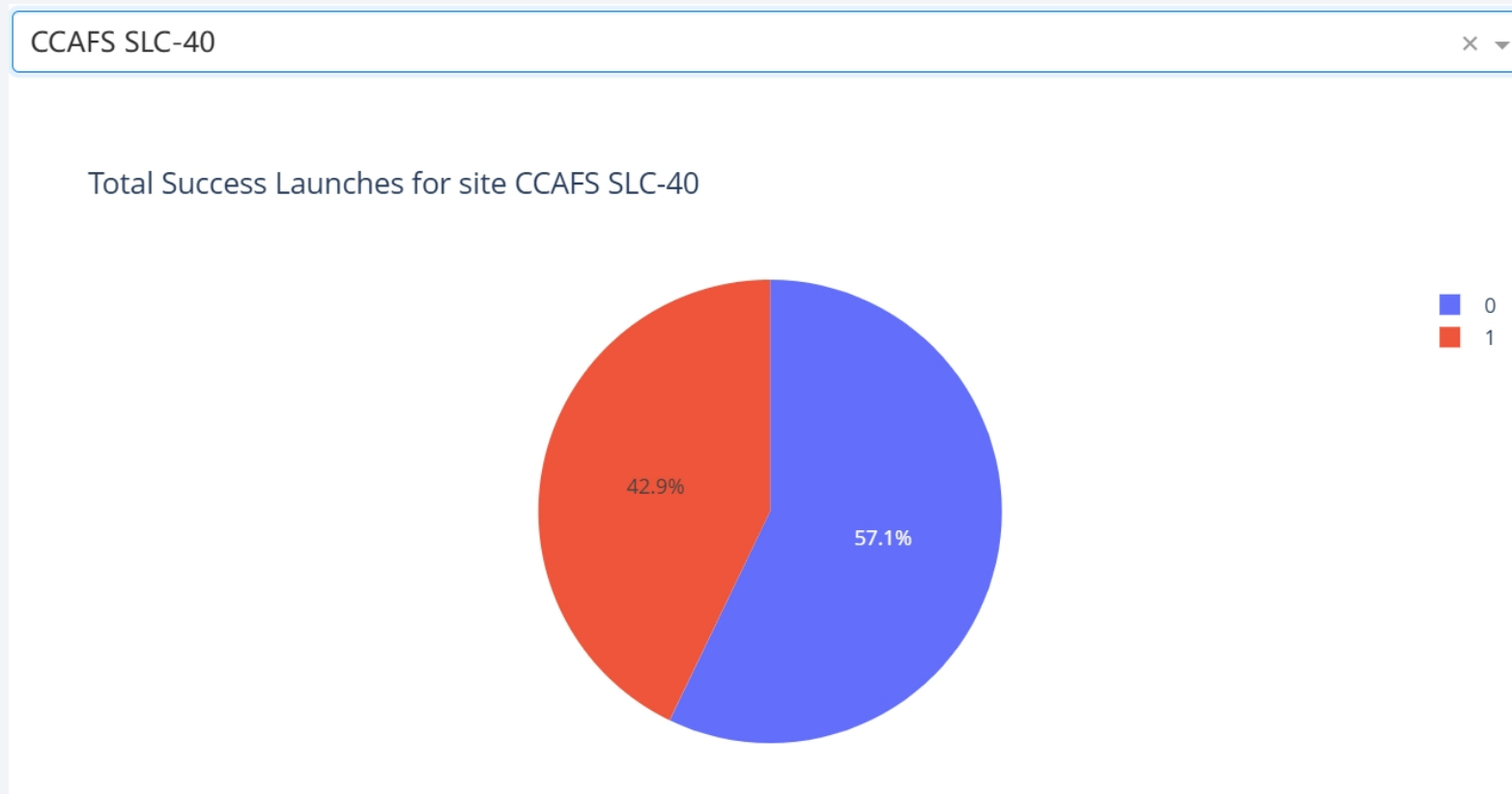
Section 4

# Build a Dashboard
# with Plotly Dash

# This is the SpaceX Dashboard

- Here you can see a dropdown option at the top to select different sites. It will change the pie chart below to your selection.

# The Launch Site with Highest Success Ratio

- In the preprocessed data, CCAFS SLC-40 has the highest success ratio

- Here, 0 represents bad outcome and 1 is the good outcome, so the success ratio of CCAFS SLC-40 is 42.9%
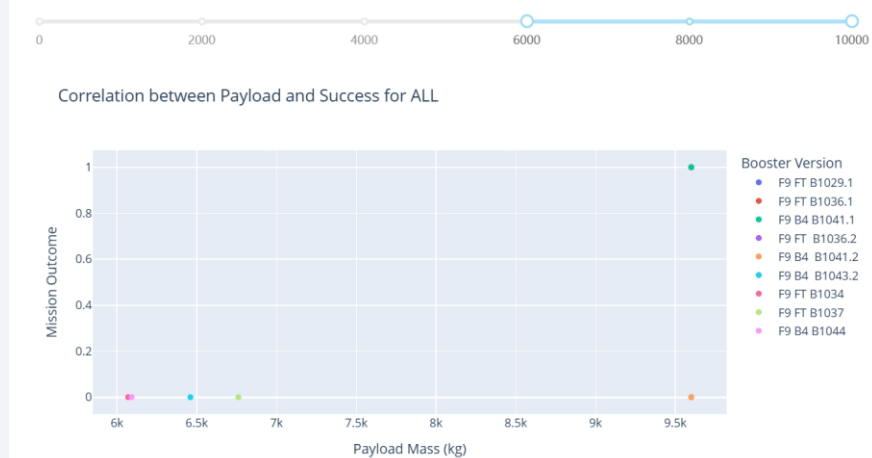
# Payload is a bigger factor



- From the scatter, we can notice that most success outcomes are from payload mass < 6000 kg range

- When payload > 6000kg, there are only 6 samples and only 1 success

Section 5

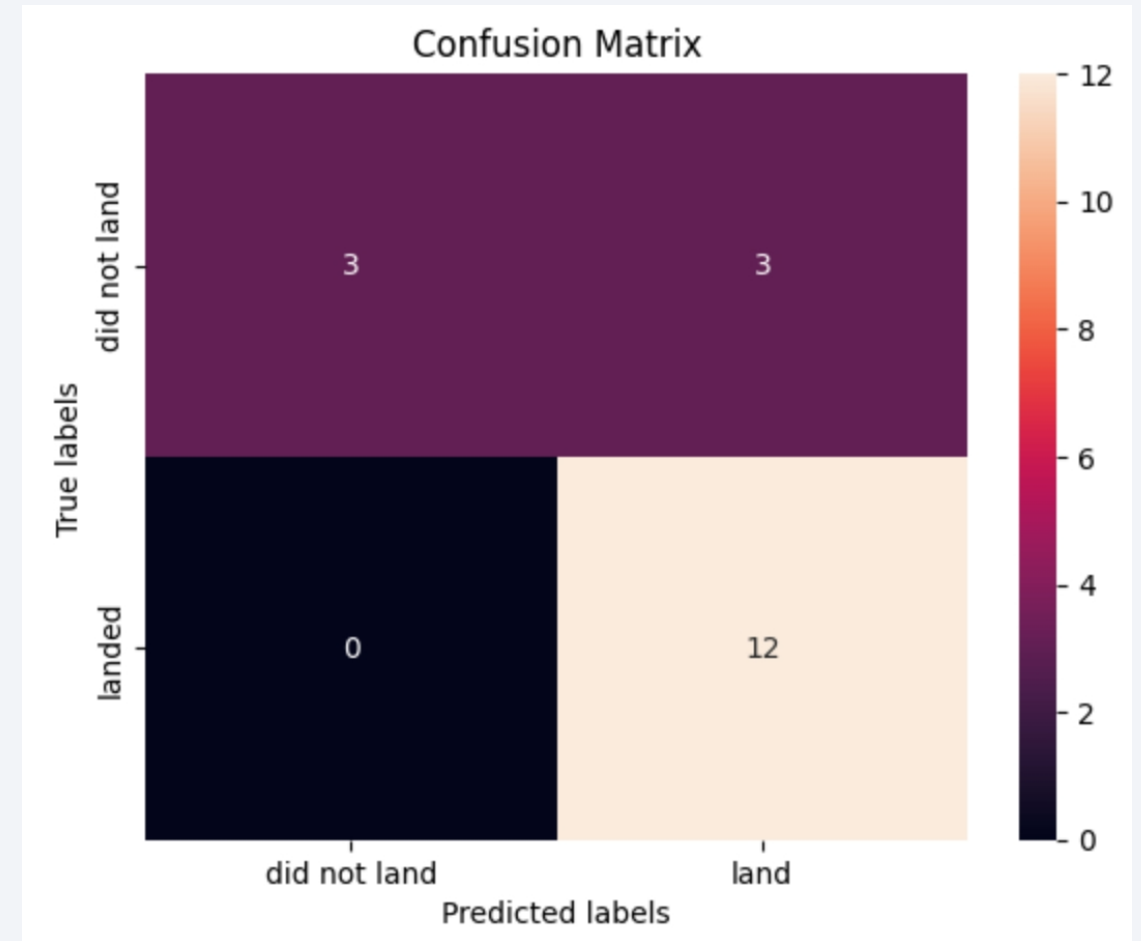# Predictive Analysis (Classification)

# Classification Accuracy

- The tree model has the highest classification accuracy for train set

- But all the models returned exactly same for test set, which means they works similar

# Confusion Matrix

- All 4 models returned the same matrix it is shown on the right

- This result works perfectly on landed group but false land still happens which may get people over confident

# Conclusions

- Through the data, we noticed that some factors such as payload mass, orbit, and time are all influential to the outcome

- With data science analysis, we can easily collect massive amount of record and turn it into one meaningful dataframe

- Through folium map, we can find some common attributes from the launch sites, for example they all located not far from the coast

- When we are trying to predict the outcome or class, all the method including logistic regression, SVM, decision tree, KNN, are doing well

- The analysis shows good potential to predict subsequent launch of SpaceX with iteration

# Appendix

- Include any relevant assets like Python code snippets, SQL queries, charts, Notebook outputs, or data sets that you may have created during this project

Thank you!