

# Winning Space Race with Data Science

PAVITHRAN. S. N

30 DEC 2023



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

# Executive Summary

3

---

## **• Summary of methodologies:**

- Data Collection using Web scraping and REST API queries
- Data Wrangling to Classify Launches based on Success and transform data into standardized numeric form
- Exploratory Data Analysis using SQL and Visualization packages for Python
- Interactive Plotly Web App to visualize payload and success launch data at each Launch Site
- Exploring Launch Sites using interactive Folium Maps
- Predictive analysis for classification of Rocket Landing Success

## **• Summary of all results:**

- Exploratory Data Analysis Results
- Predictive Analysis Results

# Introduction

---

SpaceX 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 SpaceX 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 SpaceX for a rocket launch. Thus it is advantageous to be able to predict whether the Falcon 9 first stage will land successfully for new missions.

To make valuable predictions we must solve the following:

- What factors of a mission influence Falcon 9 launch success?
- What conditions must be met by SpaceX to ensure the highest probability of success for a given mission?

Section 1

# Methodology

# Methodology

6

---

## Executive Summary

- Data collection methodology:
  - Using Space X API and web scraping from wikipedia.
- Perform data wrangling
  - Clean the data
- Perform exploratory data analysis (EDA) using visualization and SQL
- Perform interactive visual analytics using Folium and Plotly Dash
- Perform predictive analysis using classification models
  - Creating various Machine learning models and validate by score.

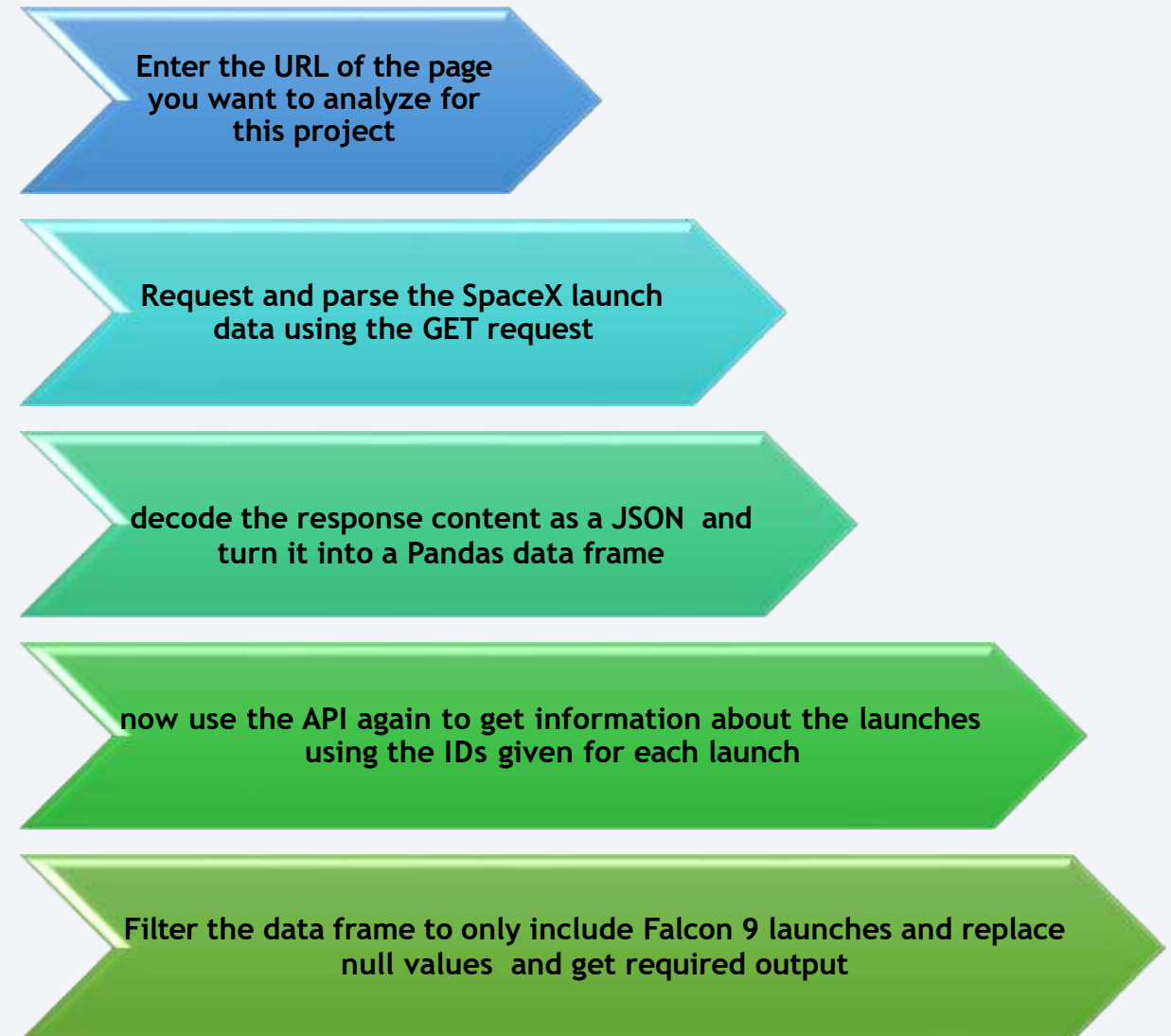


# Data Collection

7

The Data sets are collected by

- SpaceX API request.
- Web Scraping

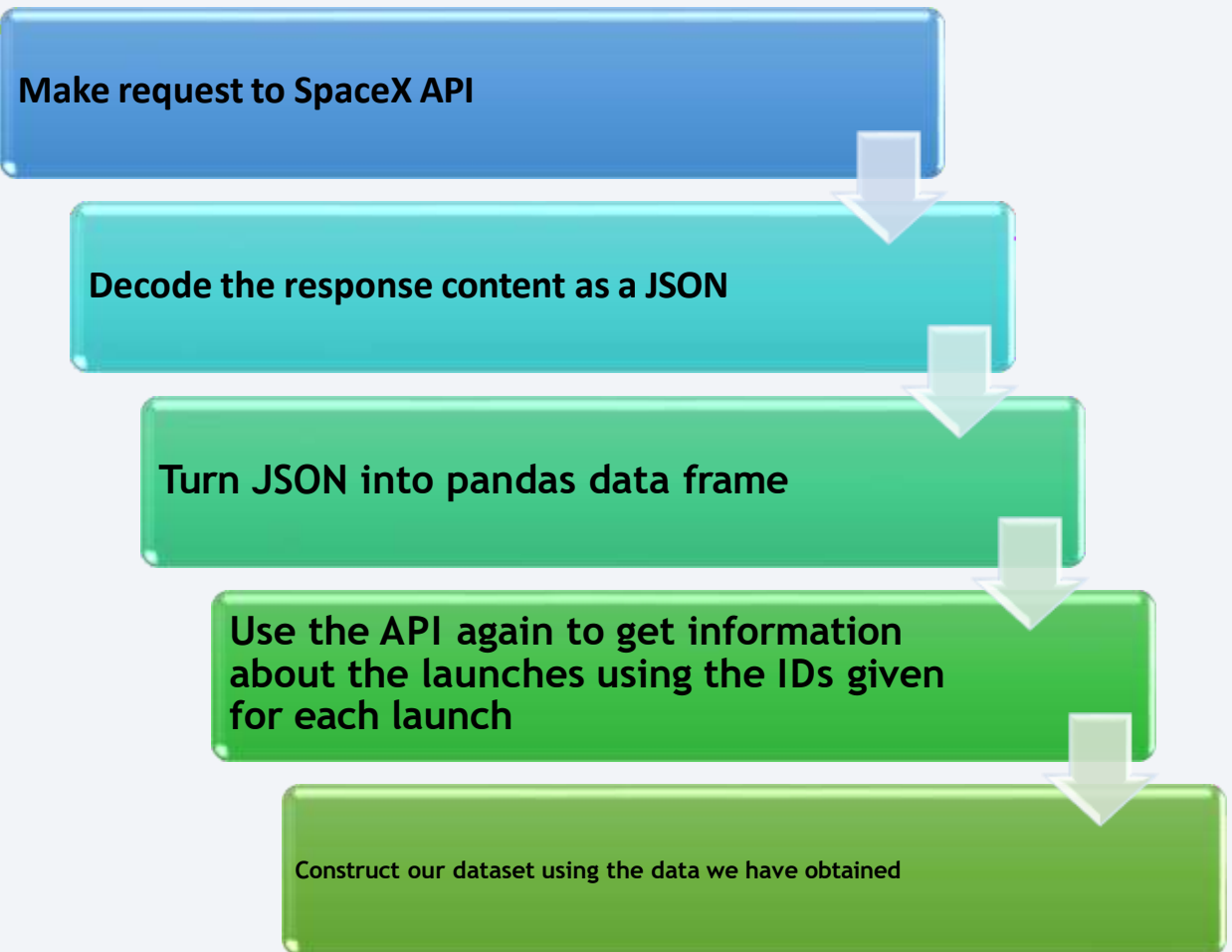


# Data Collection – SpaceX API

8

## 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



[https://github.com/imyounghman/Applied-Data-Science-Capstone-Falcon9-Spacex/blob/2034e694be52116d3a04652a0f0dd0c3ddcd94f2/1\\_jupyter-labs-spacex-data-collection-api.ipynb](https://github.com/imyounghman/Applied-Data-Science-Capstone-Falcon9-Spacex/blob/2034e694be52116d3a04652a0f0dd0c3ddcd94f2/1_jupyter-labs-spacex-data-collection-api.ipynb)

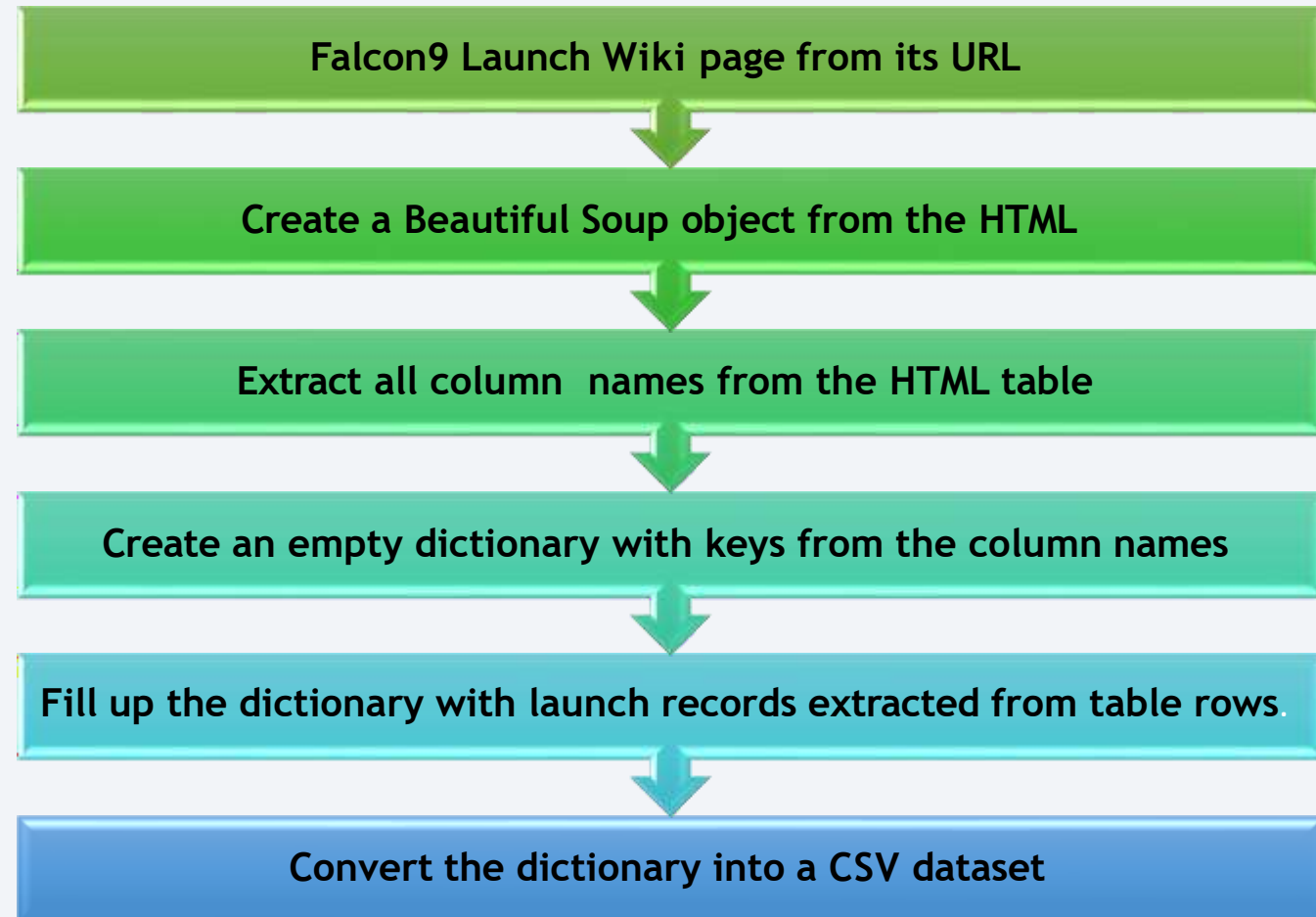


# Data Collection - Web Scrapping

9

- **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



[https://github.com/imyounghman/Applied-Data-Science-Capstone-Falcon9-Spacex/blob/2034e694be52116d3a04652a0f0dd0c3ddcd94f2/2\\_jupyter-labs-web scraping.ipynb](https://github.com/imyounghman/Applied-Data-Science-Capstone-Falcon9-Spacex/blob/2034e694be52116d3a04652a0f0dd0c3ddcd94f2/2_jupyter-labs-web scraping.ipynb)

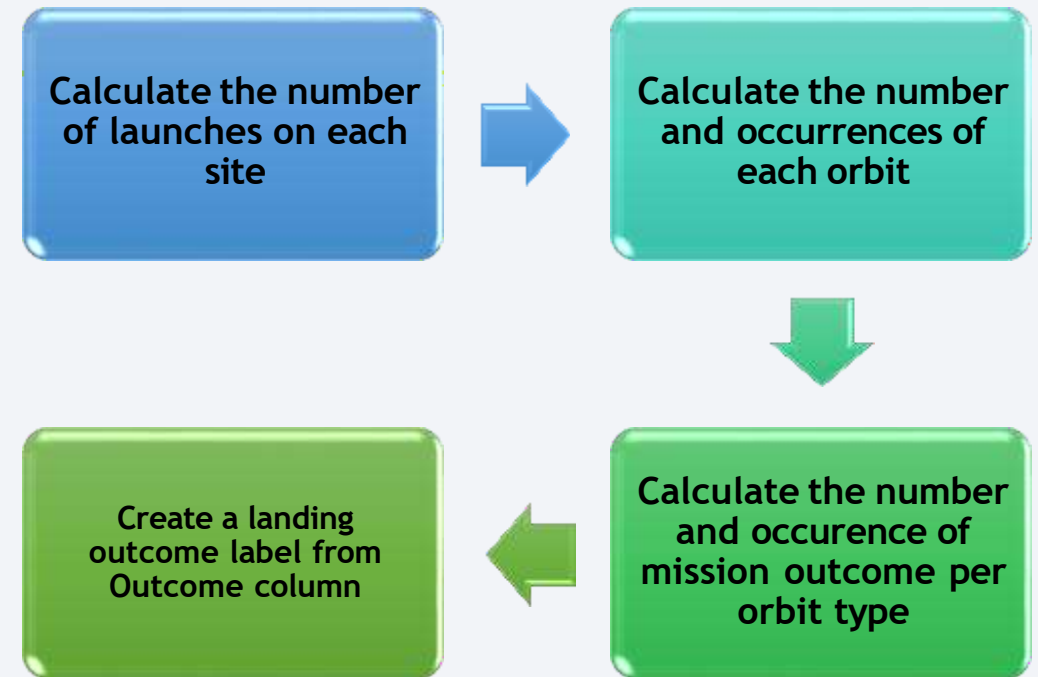
# Data Wrangling

10

➤ Data Wrangling process is given in a flow chart for a over view.

➤ **steps**

- Perform EDA and determine
- data labels
- Calculate:
  - # of launches for each site
  - # and occurrence of orbit
  - # and occurrence of mission outcome per orbit type]
- Create binary landing outcome column (dependent variable)
- Export data to csv file



[https://github.com/imyounghman/Applied-Data-Science-Capstone-Falcon9-Spacex/blob/2034e694be52116d3a04652a0f0dd0c3ddcd94f2/3\\_labs-jupyter-spacex-Data](https://github.com/imyounghman/Applied-Data-Science-Capstone-Falcon9-Spacex/blob/2034e694be52116d3a04652a0f0dd0c3ddcd94f2/3_labs-jupyter-spacex-Data)

# EDA with Data Visualization

11

## Types of Charts Used :

- **scatter plot** - Flight Number vs Payload Mass , Flight Number vs Launch Sites , Payload and Launch Sites , Flight Number and Orbit Type , Payload and Orbit Type
- **Bar chart** – Success rate of each orbit
- **Line plot** – success rate and Date

[https://github.com/imyoungman/Applied-Data-Science-Capstone-Falcon9-SpaceX/blob/2034e694be52116d3a04652a0f0dd0c3ddcd94f2/5\\_jupyter-labs-eda-dataviz.ipynb.jupyterlite.ipynb](https://github.com/imyoungman/Applied-Data-Science-Capstone-Falcon9-SpaceX/blob/2034e694be52116d3a04652a0f0dd0c3ddcd94f2/5_jupyter-labs-eda-dataviz.ipynb.jupyterlite.ipynb)

# EDA with SQL

12

---

## Summary of SQL queries that were used:

- Display the names of the unique launch sites in the space mission
- Display 5 records where launch sites begin with the string 'CCA
- Display the total payload mass carried by boosters launched by NASA (CRS)
- Display average payload mass carried by booster version F9 v1.1
- List the date when the first successful landing outcome in ground pad was achieved
- List the names of the boosters which have success in drone ship and have payload mass greater than 4000 but less than 6000
- List the total number of successful and failure mission outcomes
- List the names of the booster versions which have carried the maximum payload mass. Use a subquery
- List the failed landing outcomes in drone ship, their booster versions, and launch site names for in year 2015
- Rank the count of landing outcomes (such as Failure (drone ship) or Success (ground pad)) between the date 2010-06-04 and 2017-03-20, in descending order

[https://github.com/imyoungman/Applied-Data-Science-Capstone-Falcon9-SpaceX/blob/2034e694be52116d3a04652a0f0dd0c3ddcd94f2/4\\_jupyter-labs-eda-sql-](https://github.com/imyoungman/Applied-Data-Science-Capstone-Falcon9-SpaceX/blob/2034e694be52116d3a04652a0f0dd0c3ddcd94f2/4_jupyter-labs-eda-sql-)

# Build an Interactive Map with Folium

13

- Folium Markers were used to show the Space X launch sites and their nearest important landmarks like railways, highways, cities and coastlines.
- Polylines were used to connect the launch sites to their nearest land marks.
- **Red** represents rocket launch Failures
- **Green** represents the successes.

[https://github.com/imyoungman/Applied-Data-Science-Capstone-Falcon9-Spacex/blob/2034e694be52116d3a04652a0f0dd0c3ddcd94f2/6\\_lab\\_jupyter\\_launch\\_site\\_location.jupyterlite.ipynb](https://github.com/imyoungman/Applied-Data-Science-Capstone-Falcon9-Spacex/blob/2034e694be52116d3a04652a0f0dd0c3ddcd94f2/6_lab_jupyter_launch_site_location.jupyterlite.ipynb)

# Build a Dashboard with Plotly Dash

14

- Pie charts and scatter charts were used to visualize the launch records of Space X.
- These charts displayed the rocket launch success rate per launch site. We were able to get an understanding of the factors that may have been influencing the success rate at each site. Such as the payload mass and booster versions.
- Successful launches were represented by 1 while failures were represented by 0.

[https://github.com/imyoungman/Applied-Data-Science-Capstone-Falcon9-SpaceX/blob/2034e694be52116d3a04652a0f0dd0c3ddcd94f2/7\\_spacex\\_dash\\_app.py](https://github.com/imyoungman/Applied-Data-Science-Capstone-Falcon9-SpaceX/blob/2034e694be52116d3a04652a0f0dd0c3ddcd94f2/7_spacex_dash_app.py)

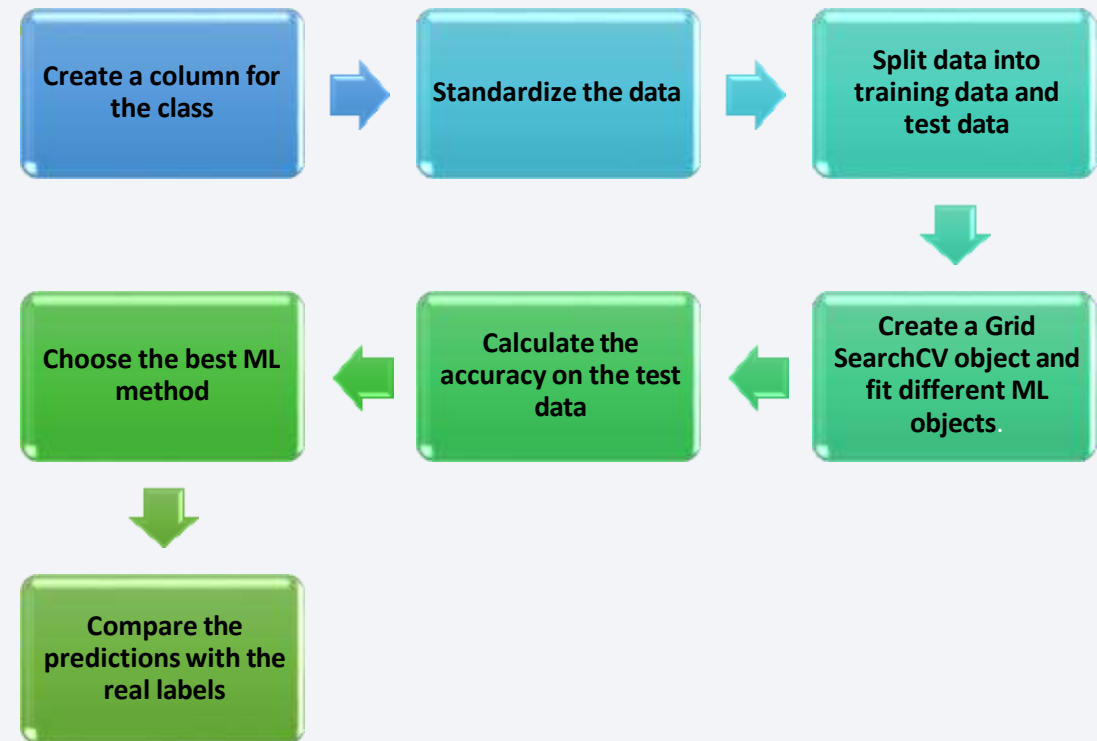


# Predictive Analysis (Classification)

15

Scikit-learn is Machine Learning library that was used for predictive analysis. The following took place:

- Created a machine learning pipeline to predict if the first stage will land given the data.



[https://github.com/imyongman/Applied-Data-Science-Capstone-Falcon9-Spacex/blob/2034e694be52116d3a04652a0f0dd0c3ddcd94f2/8\\_SpaceX\\_Machine\\_Learning\\_Prediction\\_Part\\_5.jupyterlite.ipynb](https://github.com/imyongman/Applied-Data-Science-Capstone-Falcon9-Spacex/blob/2034e694be52116d3a04652a0f0dd0c3ddcd94f2/8_SpaceX_Machine_Learning_Prediction_Part_5.jupyterlite.ipynb)

- The exploratory data analysis has shown us that successful landing outcomes are somewhat correlated with flight number. It was also apparent that successful landing outcomes have had a significant increase since the year 2015.
- All launch sites are located near the coast line. Perhaps, this makes it easier to test rocket landings in the water.
- sites are also located near highways and railways. This may facilitate transportation of equipment and research material.
- The machine learning were able to predict the landing success of rockets with an accuracy score of 83.33%.

The background of the slide is an abstract composition of numerous thin, overlapping lines and streaks in shades of blue, red, and teal. These lines are oriented diagonally, creating a sense of motion and depth. The overall effect is reminiscent of a digital data visualization or a stylized representation of a complex system.

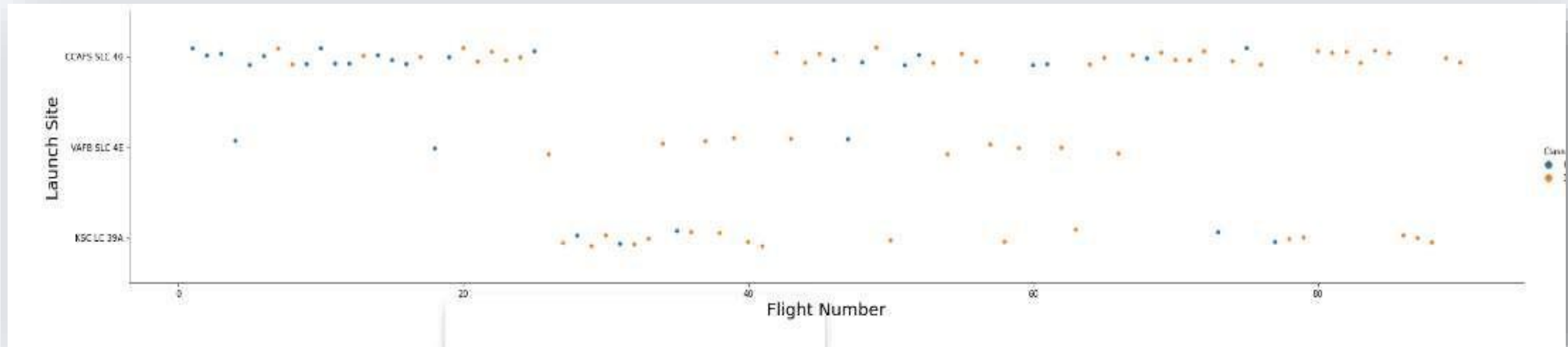
Section 2

# Insights drawn from EDA



# Flight Number vs. Launch Site

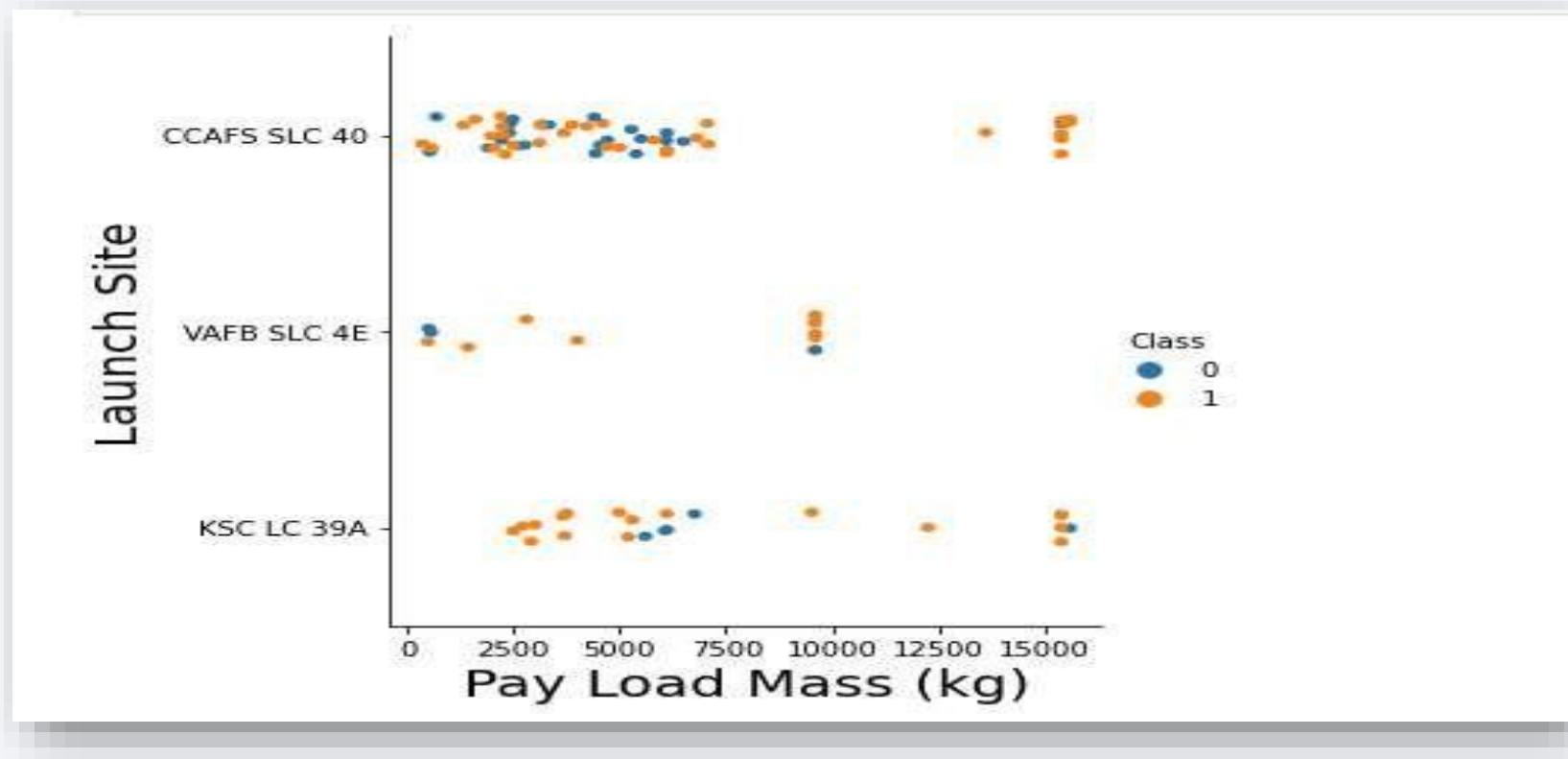
18



- It appears that there were more successful landings as the flight numbers increased. launch site **CCAFS SLC 40** had the most number of landing.

# Payload vs. Launch Site

19



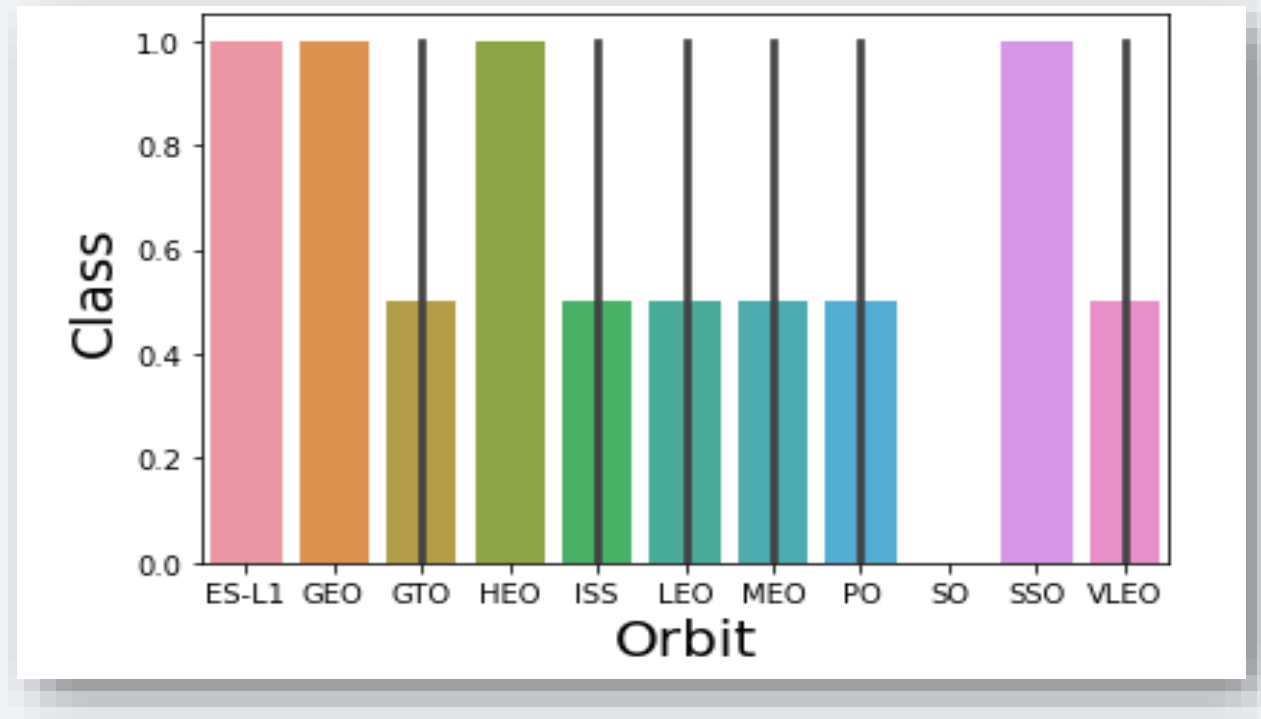
- Now if you observe the scatter point chart, you will find for the VAFB-SLC launch site there are no rockets launched for heavy payload mass(greater than 10000).

# Success Rate vs. Orbit Type

20

The highest success rate ORBITS are

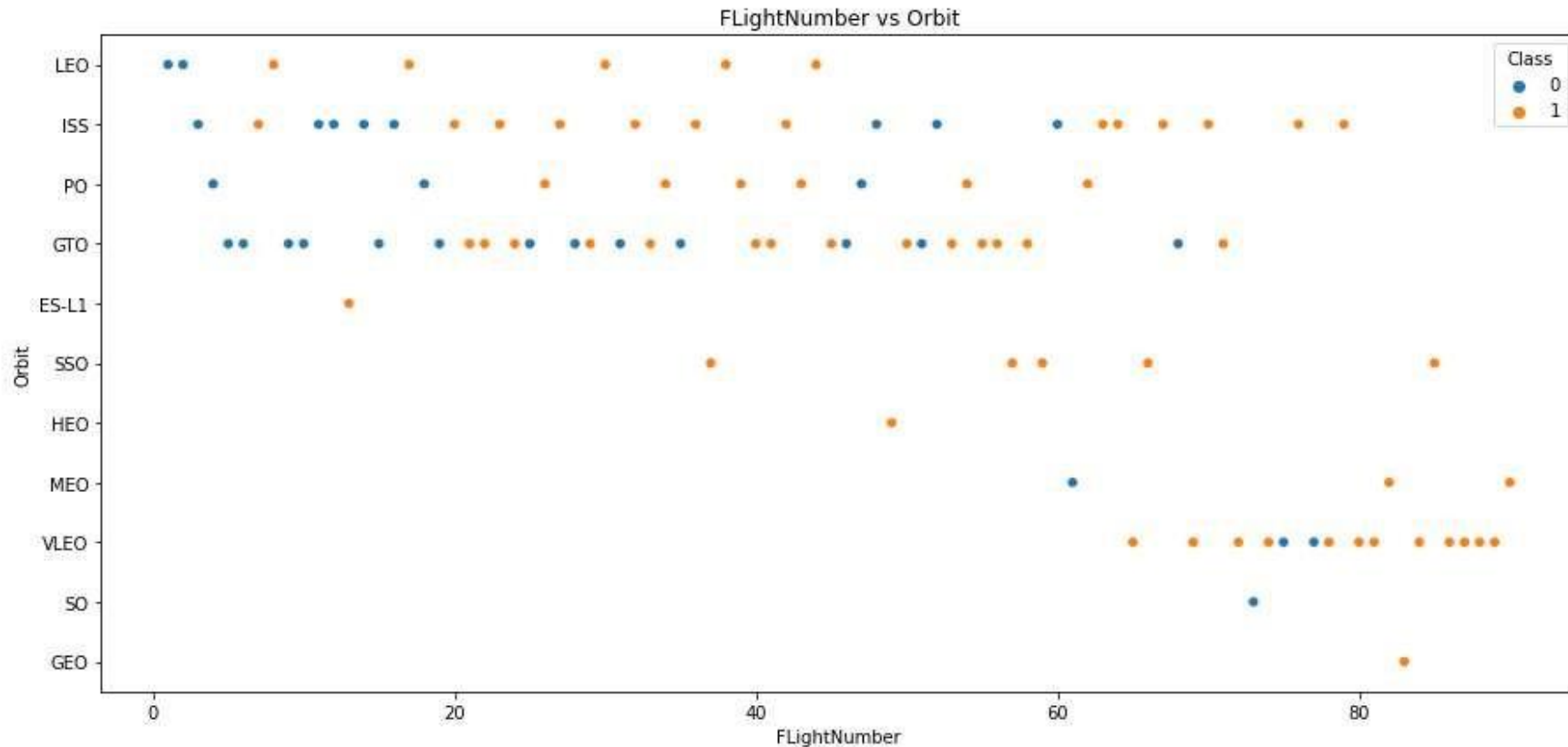
**ES-L1 GEO SSO HEO**





# Flight Number vs. Orbit Type

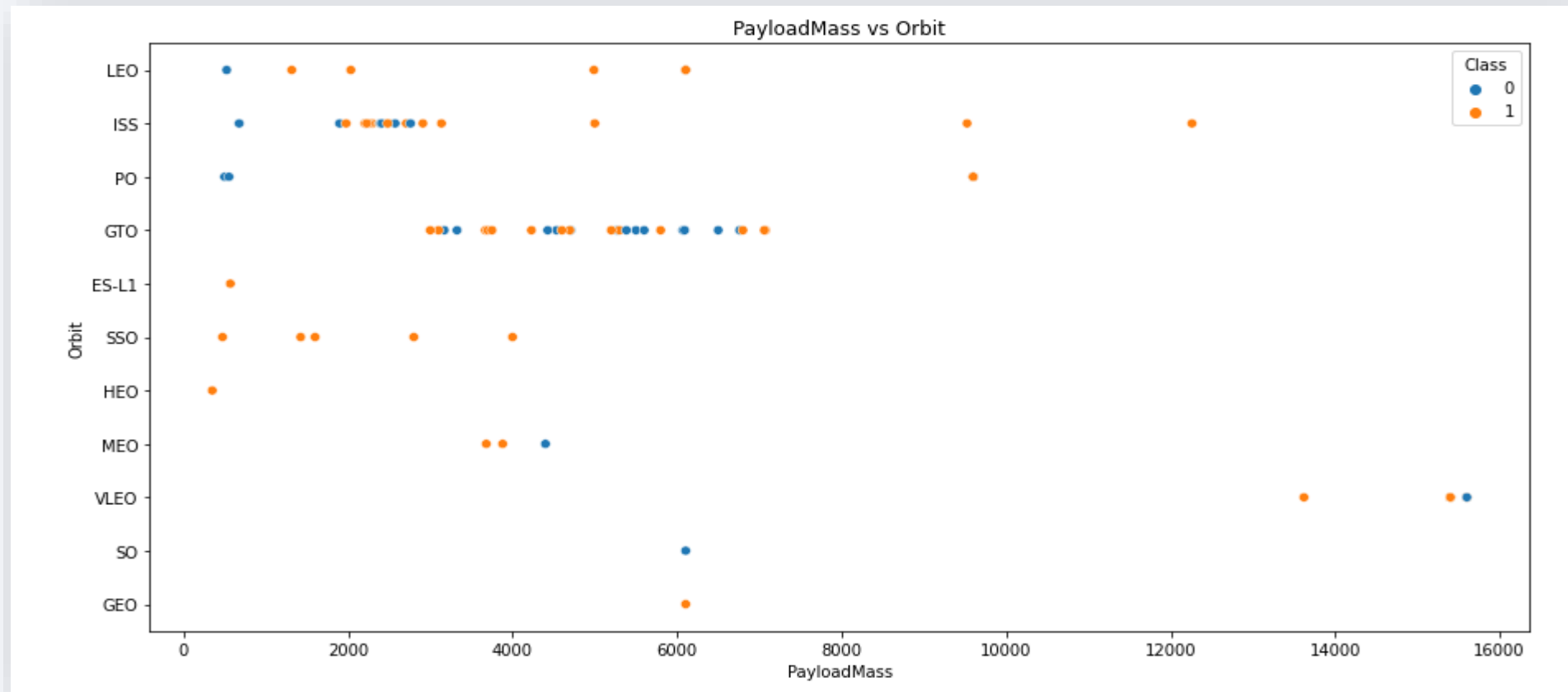
21



You can see that in the LEO orbit the Success appears related to the number of flights; on the other hand, there seems to be no relationship between flight number when in GTO orbit.

# Payload vs. Orbit Type

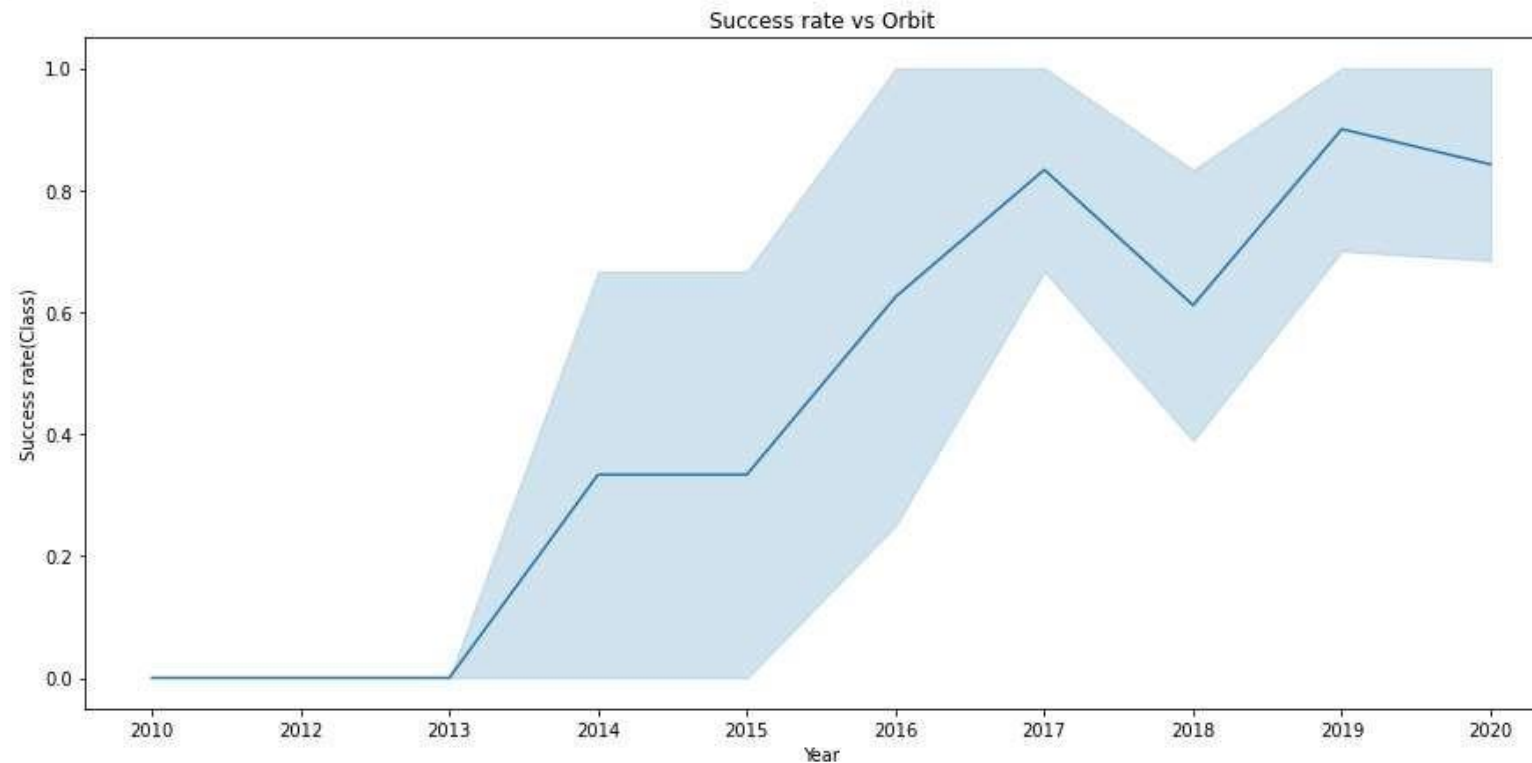
22



- With heavy payloads the successful landing or positive landing rate are more for Polar, LEO and ISS.
- However for GTO we cannot distinguish this well as both positive landing rate and negative landing(unsuccesful mission) are both there.

# Launch Success Yearly Trend

23



It is apparent that the success rate has significantly increased from 2013 to 2020.

# All Launch Site Names

---

24

Given the data, these are the names of the launch sites where different rocket landings were attempted:

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

# Launch Site Names Beginning with 'CCA'

25

```
In [18]: %sql SELECT * FROM SPACEXTBL WHERE LAUNCH_SITE LIKE 'CCA%' LIMIT 5;
```

```
* ibm_db_sa://gfd86828:***@3883e7e4-18f5-4afe-be8c-fa31c41761d2.bs2io90108kqb1od81cg.databases.appdomain.cloud:31498/bludb
Done.
```

```
Out[18]:
```

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	07:44:00	F9 v1.0 B0005	CCAFS LC-40	Dragon demo flight C2	525	LEO (ISS)	NASA (COTS)	Success	No attempt
2012-10-08	00:35:00	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

These are 5 records where launch sites begin with the letters 'CCA'. As we can see, there are other organizations besides Space X that were testing their rockets.

# Total Payload Mass

26

Display the total payload mass carried by boosters launched by NASA (CRS)

```
In [23]: %sql SELECT SUM(PAYLOAD_MASS__KG_) FROM SPACEXTBL WHERE CUSTOMER = 'NASA(CRS)';  
  
* ibm_db_sa://gfd86828:***@3883e7e4-18f5-4afe-be8c-fa31c41761d2.bs2io90108kqb1od81cg.databases.appdomain.cloud:31498/bludb  
Done.  
Out[23]: 1
```

- The information in the picture displays the total payload mass carried by boosters launched by NASA



# Average Payload Mass by F9 v1.1

27

Display average payload mass carried by booster version F9 v1.1

```
In [24]: %sql SELECT AVG(PAYLOAD_MASS_KG_) FROM SPACEXTBL WHERE BOOSTER_VERSION = 'F9 v1.1'
* ibm_db_sa://gfd86828:***@3883e7e4-18f5-4afe-be8c-fa31c41761d2.bs2io90108kqb1od81cg.databases.appdomain.cloud:31498/bludb
Done.
Out[24]: 1
2928
```

- The average payload mass carried by F9 v1.1 was 2928.4 kg.

# First Successful Ground Landing Date

28

List the date when the first successful landing outcome in ground pad was achieved.

*Hint: Use min function*

```
In [28]: %sql select min(DATE) from SPACEXTBL where Landing__Outcome = 'Success (ground pad)';

* ibm_db_sa://gfd86828:***@3883e7e4-18f5-4afe-be8c-fa31c41761d2.bs2io90l08kqb1od8l1cg.databases.appdomain.cloud:31498/bludb
Done.

Out[28]:      1
          2015-12-22
```

- From the picture given above you can see that the first successful ground pad was in 22 December 2015.

# Successful Drone Ship Landing with Payload between 4000 and 6000

29

List the names of the boosters which have success in drone ship and have payload mass greater than 4000 but less than 6000

```
In [32]: %sql SELECT BOOSTER_VERSION from SPACEXTBL WHERE LANDING__OUTCOME = 'Success (drone ship)' and PAYLOAD_MASS__KG_ >4000 and PAYLOAD_MASS__KG_ <6000;

* ibm_db_sa://gfd86828:***@3883e7e4-18f5-4afe-be8c-fa31c41761d2.bs2io90l08kqb1od8l1cg.databases.appdomain.cloud:31498/bludb
Done.

Out[32]: booster_version
         F9 FT B1022
         F9 FT B1026
         F9 FT B1021.2
         F9 FT B1031.2
```

- It appears that there only 4 Boosters with a payload mass between 4000 and 6000 they are
- F9 FT B1022
- F9 FT B1026
- F9 FT B1021.2
- F9 FT B1031.2

# Total Number of Successful and Failure Mission Outcomes

30

List the total number of successful and failure mission outcomes

```
In [33]: %sql select count(MISSION_OUTCOME) from SPACEXTBL where MISSION_OUTCOME = 'Success' or MISSION_OUTCOME = 'Failure (in flight)'
```

```
* ibm_db_sa://gfd86828:***@3883e7e4-18f5-4afe-be8c-fa31c41761d2.bs2io90l08kqb1od8l1cg.databases.appdomain.cloud:31498/bludb  
Done.
```

```
Out[33]: 1
```

```
100
```

- The Above picture show the total number of successful and failure mission outcomes

# Boosters That Carried the Maximum Payload Mass

31

List the names of the booster\_versions which have carried the maximum payload mass. Use a subquery

```
In [34]: %sql SELECT BOOSTER_VERSION FROM SPACEXTBL WHERE PAYLOAD_MASS__KG_ = (SELECT max(PAYLOAD_MASS__KG_) FROM SPACEXTBL);
```

\* ibm\_db\_sa://gfd86828:\*\*\*@3883e7e4-18f5-4afe-be8c-fa31c41761d2.bs2io90l08kqb1od8lcg.databases.appdomain.cloud:31498/bludb  
Done.

Out[34]: **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

- From the above picture it shows that 12 boosters have carried the maximum payload mass of 15600 kg.

# 2015 Launch Records - Failed Landing Outcomes

32

List the failed landing\_outcomes in drone ship, their booster versions, and launch site names for in year 2015

```
In [18]: task_9 = '''
          SELECT BoosterVersion, LaunchSite, LandingOutcome
          FROM SpaceX
          WHERE LandingOutcome LIKE 'Failure (drone ship)'
          AND Date BETWEEN '2015-01-01' AND '2015-12-31'
          ...
          create_pandas_df(task_9, database=conn)
```

```
Out[18]:
```

	boosterversion	launchsite	landingoutcome
0	F9 v1.1 B1012	CCAFS LC-40	Failure (drone ship)
1	F9 v1.1 B1015	CCAFS LC-40	Failure (drone ship)

- 2 boosters **F9 v1.1B1012\_CCAFS LC-40** and **F9v1.1B1015 CCAFS LC-40** failed to land at 2015



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

33

Rank the count of landing outcomes (such as Failure (drone ship) or Success (ground pad)) between the date 2010-06-04 and 2017-03-20, in descending order

```
In [42]: %sql select * from SPACEXTBL where Landing__Outcome = 'Success (ground pad)' or and (DATE between '2010-06-04' and '2017-03-20') order by date desc
```

```
* ibm_db_sa://gfd86828:***@3883e7e4-18f5-4afe-be8c-fa31c41761d2.bs2io90l08kqb1od8l1cg.databases.appdomain.cloud:31498/bludb
Done.
```

```
Out[42]:
```

DATE	time_utc	booster_version	launch_site	payload	payload_mass_kg	orbit	customer	mission_outcome	landing_outcome
2017-02-19	14:39:00	F9 FT B1031.1	KSC LC-39A	SpaceX CRS-10	2490	LEO (ISS)	NASA (CRS)	Success	Success (ground pad)
2016-07-18	04:45:00	F9 FT B1025.1	CCAFS LC-40	SpaceX CRS-9	2257	LEO (ISS)	NASA (CRS)	Success	Success (ground pad)
2015-12-22	01:29:00	F9 FT B1019	CCAFS LC-40	OG2 Mission 2 11 Orbcomm-OG2 satellites	2034	LEO	Orbcomm	Success	Success (ground pad)

- The number of successful landings have increased since 2015.

A satellite view of Earth from space, showing the curvature of the planet and city lights at night. The image is a composite of a solid blue gradient on the left and a satellite photograph of Earth on the right. The Earth's surface is dark, with numerous bright yellow and orange lights representing city lights at night. The horizon line of the Earth is visible, separating the dark surface from the deep blue of the sky.

Section 3

# Launch Sites Proximities Analysis

# Launch Site Locations

35

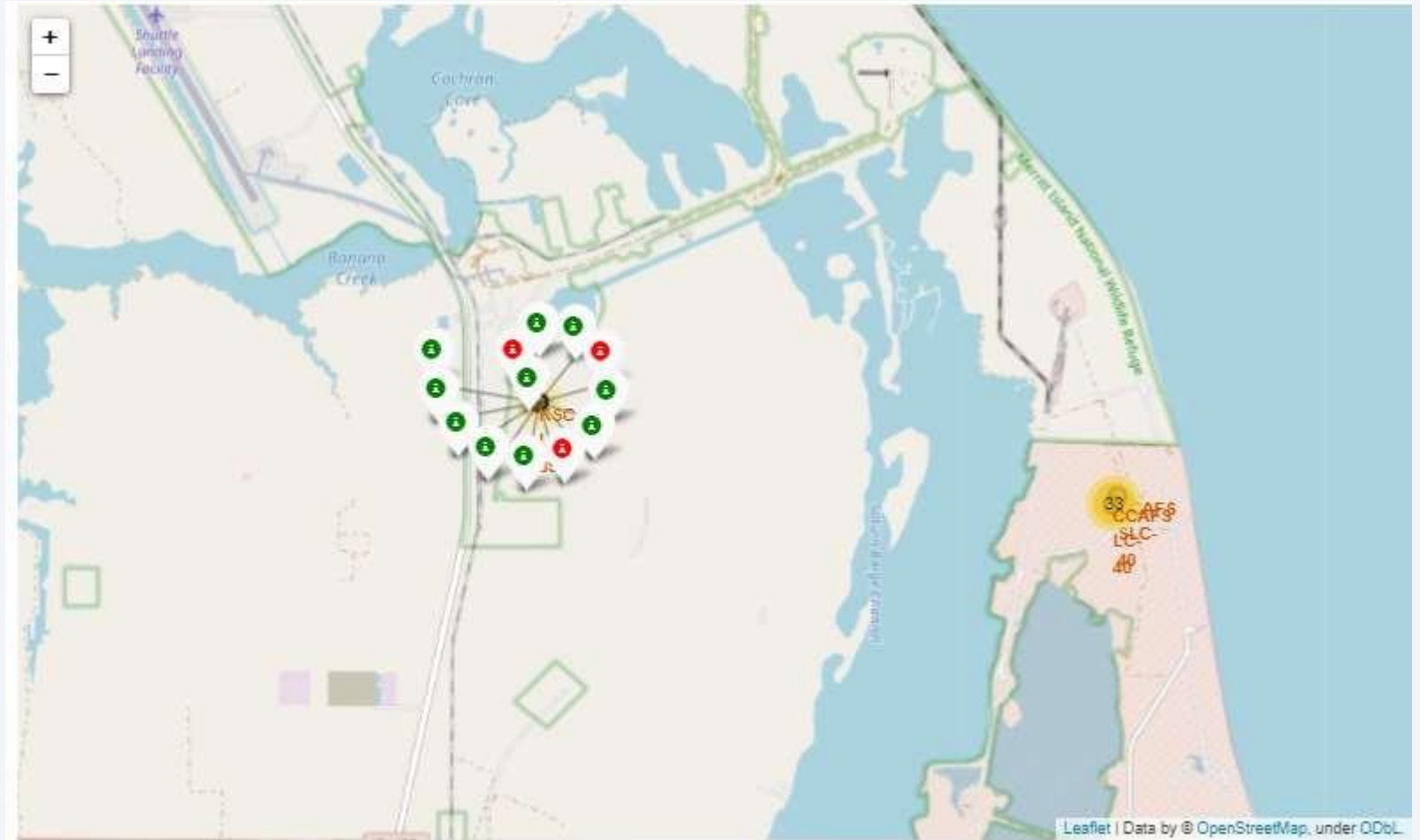
- all launch sites are in very close proximity to the coast and they are also a couple thousand kilometers away from the equator line.



# Success Rate of Rocket Launches

36

- The successful launches are represented by a **green** marker while the **red** marker represents failed rocket launches.





# Surrounding Landmarks

37

- It appears that launch sites are usually set up at least 18 km away from cities. This may be because of the desire to prevent any crashes near populated areas.
- It is also apparent that launch sites are in very close proximity to railways and highways. Perhaps, due to the necessary transportation requirements for rocket parts.
- The sites are close the coast line. This is evident with the many rocket landing tests on water bodies like the ocean.



Map Object	Colour
Nearest Highway	Green
Nearest Railway	Purple
Nearest City	Crimson
Nearest Coastline	Dark Blue

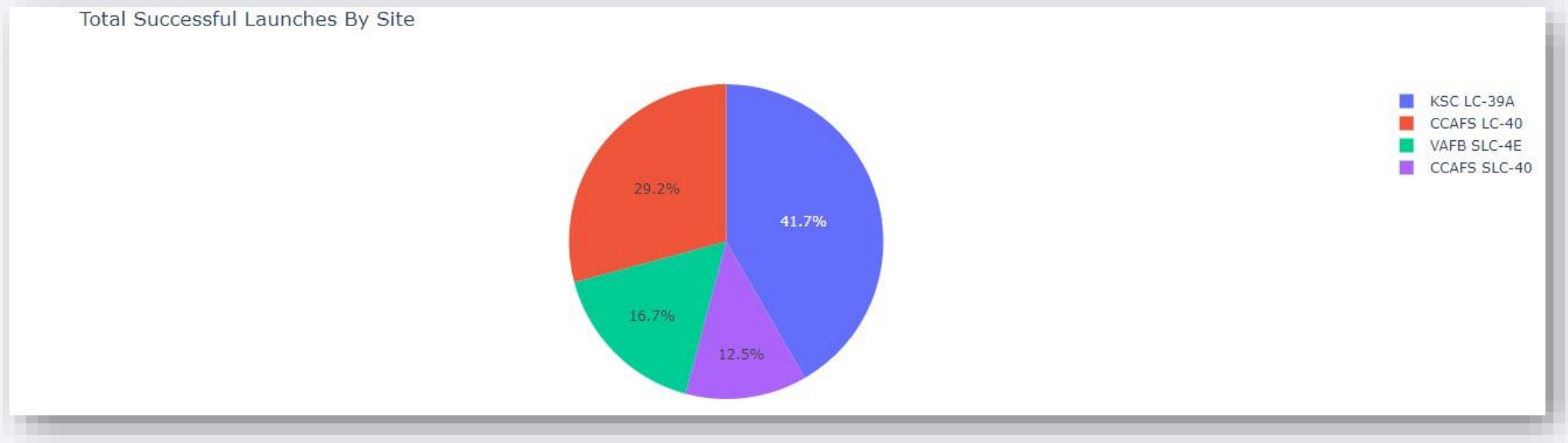


Section 4

# Build a Dashboard with Plotly Dash

# Successful Launches by Site

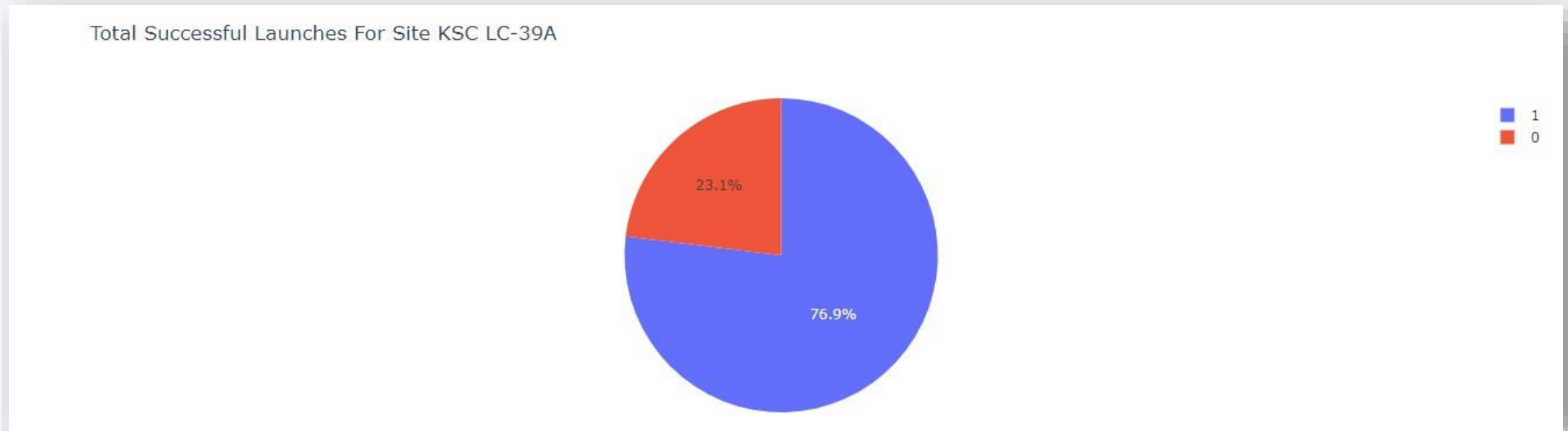
39



- You can see from the plot that Site KSC LC-39A has the largest successful launches as well the highest launch success rate.

# Total Successful Launches for Site KSC LC-39A

40

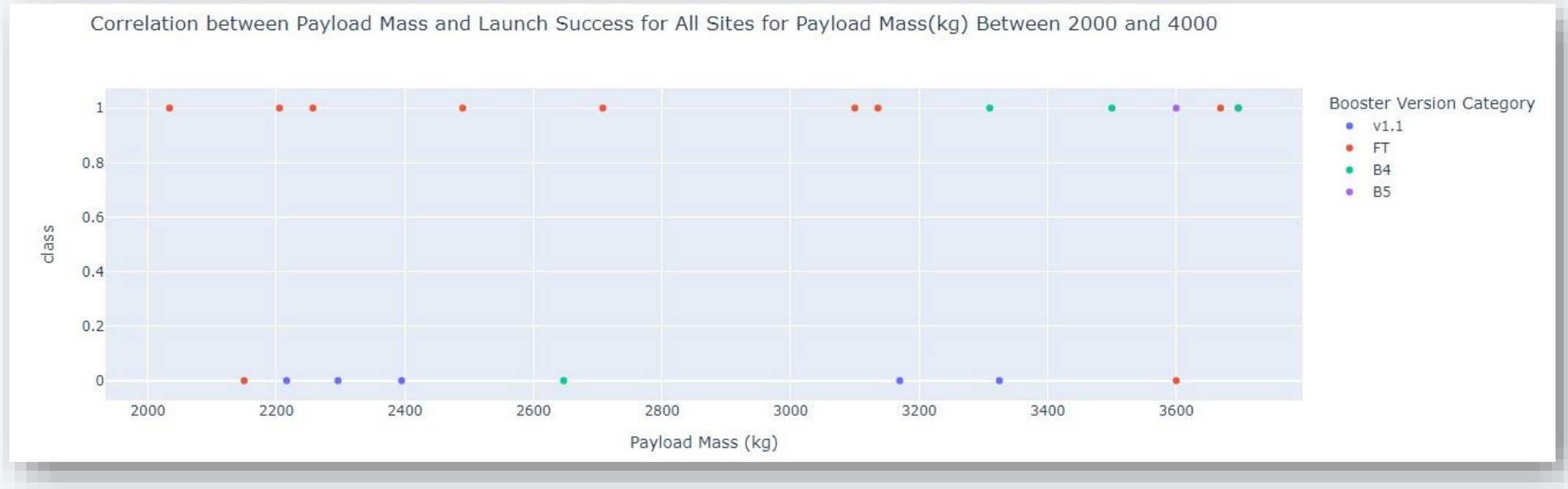


- You can see that 76.9% of the total launches at site KSC LC-39A were successful. This is the highest success rate of all the different launch sites.



# Payload Mass vs. Launch Success for All Sites

41



- It appears that the payload range between 2000 kg and 4000 kg has the highest success rate.



Section 5

# Predictive Analysis (Classification)

# Classification Accuracy

43

Find the method performs best:

```
In [28]: accuracy = [svm_cv_score, logreg_score, knn_cv_score, tree_cv_score]
accuracy = [i * 100 for i in accuracy]

method = ['Support Vector Machine', 'Logistic Regression', 'K Nearest Neighbour', 'Decision Tree']
models = {'ML Method':method, 'Accuracy Score (%)':accuracy}

ML_df = pd.DataFrame(models)
ML_df
```

```
Out[28]:
```

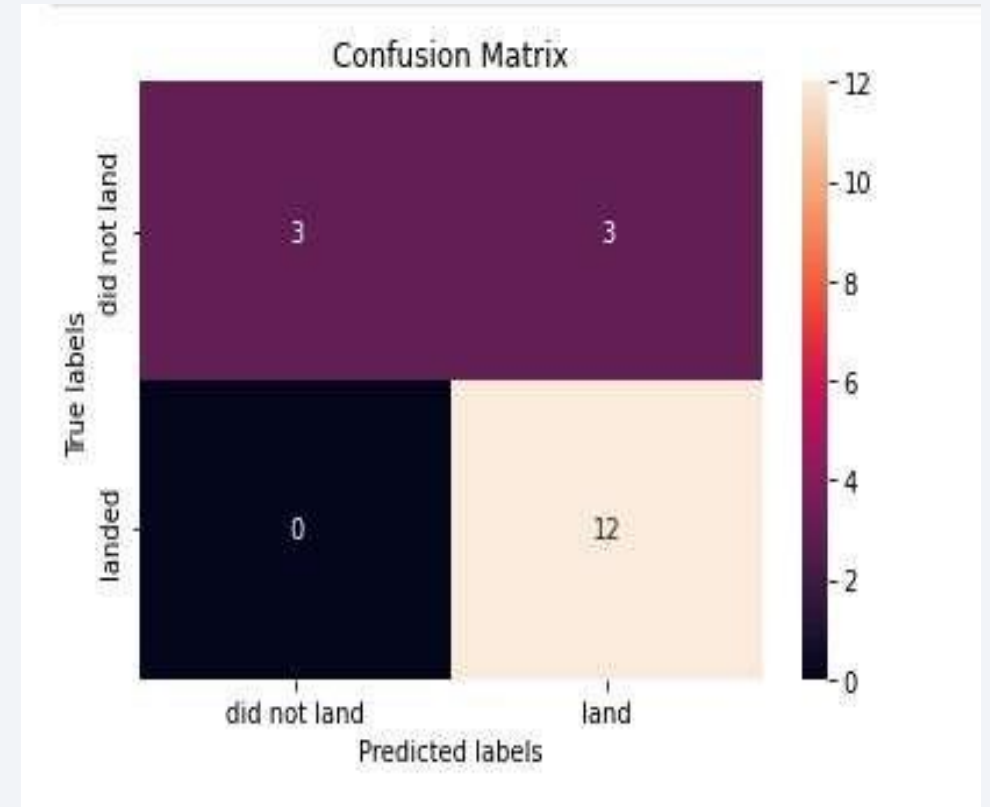
	ML Method	Accuracy Score (%)
0	Support Vector Machine	83.333333
1	Logistic Regression	83.333333
2	K Nearest Neighbour	83.333333
3	Decision Tree	83.333333

- You can see that All the methods have an identical accuracy score of 83.33%, so we decided to use Logistic Regression for the classification

# Confusion Matrix

44

- The chart shows the confusion matrix of the Logistic Regression model that was chosen.
- The model only failed to accurately predict 3 labels.



# Conclusions

---

45

In order to compete with Space X Through this process, a general picture of their success methods are

- All their launch sites are located near the coast, away from nearby cities. This enabled to them to test their rocket landings without much interference.
- Site KSC LC-39A had the highest launch success rate out of all the launch sites.
- From 2015 onwards, the success rate of rocket landings significantly increased. It was also apparent that landing success increased with flight number

All this data was used to train a machine learning model that is able to predict the landing outcome of rocket launches with 83.33% accuracy.

Haversine's Formula used to calculate distances on Folium Maps

```
from math import sin, cos, sqrt, atan2, radians

def calculate_distance(lat1, lon1, lat2, lon2):
    # approximate radius of earth in km
    R = 6373.0

    lat1 = radians(lat1)
    lon1 = radians(lon1)
    lat2 = radians(lat2)
    lon2 = radians(lon2)

    dlon = lon2 - lon1
    dlat = lat2 - lat1

    a = sin(dlat / 2)**2 + cos(lat1) * cos(lat2) * sin(dlon / 2)**2
    c = 2 * atan2(sqrt(a), sqrt(1 - a))

    distance = R * c
    return distance
```

## Things to Consider

- Dataset: A larger dataset will help build on the predictive analytics results to help understand if the findings can be generalizable to a larger data set
- Feature Analysis / PCA: Additional feature analysis or principal component analysis should be conducted to see if it can help improve accuracy
- XGBoost: Is a powerful model which was not utilized in this study. It would be interesting to see if it outperforms the other classification models

Thank you!

