



IBM Developer  
SKILLS NETWORK

# Winning Space Race with Data Science

NAME- VENKAT PRASAD DANDSENA

DATE- 09/11/2022



Executive Summary

Introduction

Methodology

Results

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**

- Machine Learning Prediction
- Exploratory Data Analysis result
- Interactive analytics in screenshots
- Predictive Analytics result



# 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.

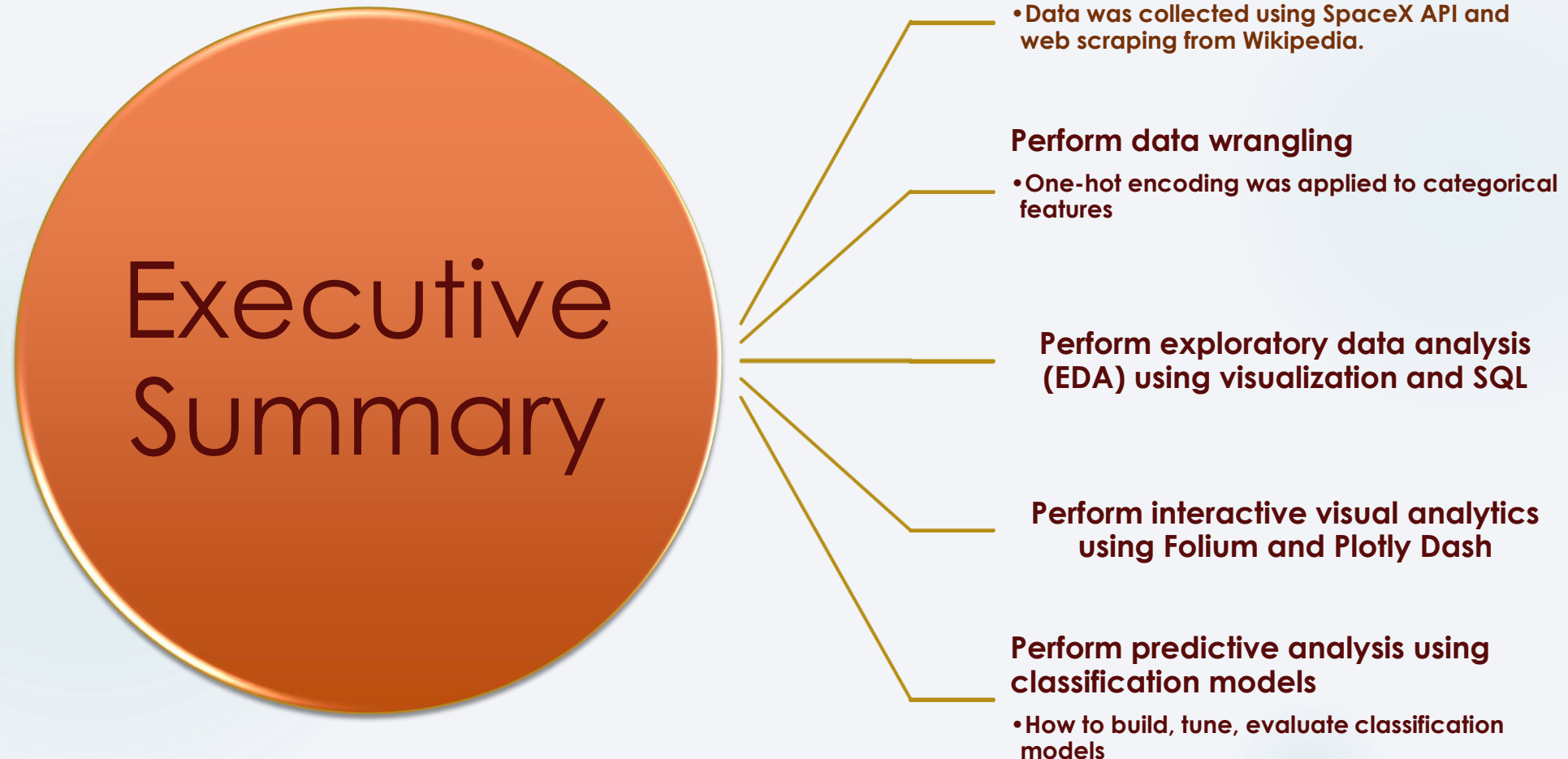
- **Problems you want to find answers**

1. What factors determine if the rocket will land successfully?
2. The interaction amongst various features that determine the success rate of a successful landing.
3. What operating conditions needs to be in place to ensure a successful landing program



Section 1

# Methodology



- 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

8

- ▶ We used the get request to the SpaceX API to collect data, clean the requested data and did some basic data wrangling and formatting.
- ▶ GitHub URL of the completed SpaceX API calls notebook ([https://github.com/Venkatdandsena/testrepo/blob/a02b100cea689267f5a055115f70292a6b51cf30/capstone\\_data\\_collectionAPI.ipynb](https://github.com/Venkatdandsena/testrepo/blob/a02b100cea689267f5a055115f70292a6b51cf30/capstone_data_collectionAPI.ipynb))

```
1. Get request for rocket launch data using API

In [6]: spacex_url="https://api.spacexdata.com/v4/launches/past"

In [7]: response = requests.get(spacex_url)

2. Use json_normalize method to convert json result to dataframe

In [12]: # Use json_normalize method to convert the json result into a dataframe
          # decode response content as json
          static_json_df = res.json()

In [13]: # apply json_normalize
          data = pd.json_normalize(static_json_df)

3. We then performed data cleaning and filling in the missing values

In [30]: rows = data_falcon9['PayloadMass'].values.tolist()[0]

          df_rows = pd.DataFrame(rows)
          df_rows = df_rows.replace(np.nan, PayloadMass)

          data_falcon9['PayloadMass'][0] = df_rows.values
          data_falcon9
```



# Data Collection - Scraping

- ▶ We applied web scrapping to webscrap Falcon 9 launch records with BeautifulSoup
- ▶ We parsed the table and converted it into a pandas dataframe
- ▶ GitHub URL of the completed web scraping notebook(<https://github.com/Venkatdandsena/testrepo/blob/a02b100cea689267f5a055115f70292a6b51cf30/Data%20Collection%20with%20Web%20Scraping%20lab.ipynb>)

```
1]: static_url = "https://en.wikipedia.org/w/index.php?title=List_of_Falcon_9_and_Falcon_Heavy_launches&oldid=1027686922"
```

Next, request the HTML page from the above URL and get a `response` object

## TASK 1: Request the Falcon9 Launch Wiki page from its URL

First, let's perform an HTTP GET method to request the Falcon9 Launch HTML page, as an HTTP response.

```
3]: # use requests.get() method with the provided static_url
response=requests.get(static_url)
# assign the response to a object
```

Create a `BeautifulSoup` object from the HTML `response`

```
5]: # Use BeautifulSoup() to create a BeautifulSoup object from a response text content
soup=BeautifulSoup(response.text, 'html')
```

Print the page title to verify if the `BeautifulSoup` object was created properly

```
7]: # Use soup.title attribute
soup.title
```

```
7]: <title>List of Falcon 9 and Falcon Heavy launches - Wikipedia</title>
```

## TASK 2: Extract all column/variable names from the HTML table header

Next, we want to collect all relevant column names from the HTML table header

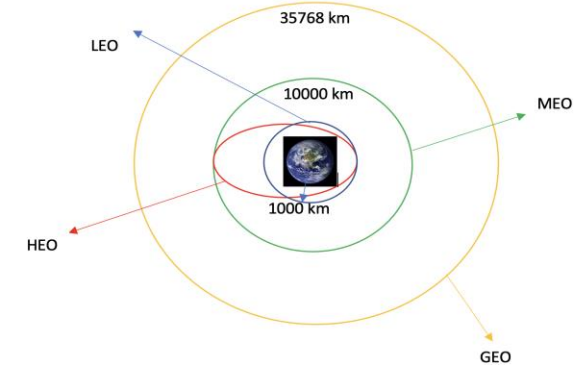
Let's try to find all tables on the wiki page first. If you need to refresh your memory about `BeautifulSoup`, please check the external reference link towards the this lab

```
3]: # Use the find_all function in the BeautifulSoup object, with element type 'table'
# Assign the result to a list called 'html_tables'
html_tables=soup.find_all('table')
```

Starting from the third table is our target table contains the actual launch records.

# 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 cs
- ▶ GitHub URL
- ▶ <https://github.com/Venkatdandsena/testrepo/blob/a02b100cea689267f5a055115f70292a6b51cf30/EDA%20lab.ipynb>



## Data Analysis

Load Space X dataset, from last section.

```
df=pd.read_csv("https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBM-DS0321EN-SkillsNetwork/datasets/dataset_part_1.csv")
df.head(10)
```

[illegible]

# EDA with Data Visualization

11

- ▶ 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

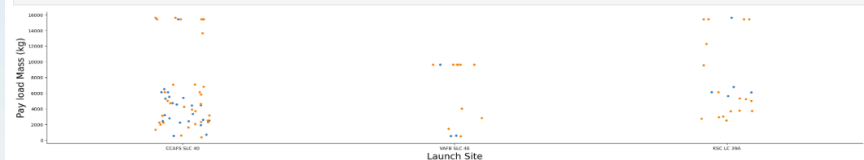
- ▶ GitHub URL

- ▶ [https://github.com/Venkat-dandsena/testrepo/blob/a02b100cea689267f5a055115f70292a6b51cf30/EDA%20visualization%20\(1\).ipynb](https://github.com/Venkat-dandsena/testrepo/blob/a02b100cea689267f5a055115f70292a6b51cf30/EDA%20visualization%20(1).ipynb)

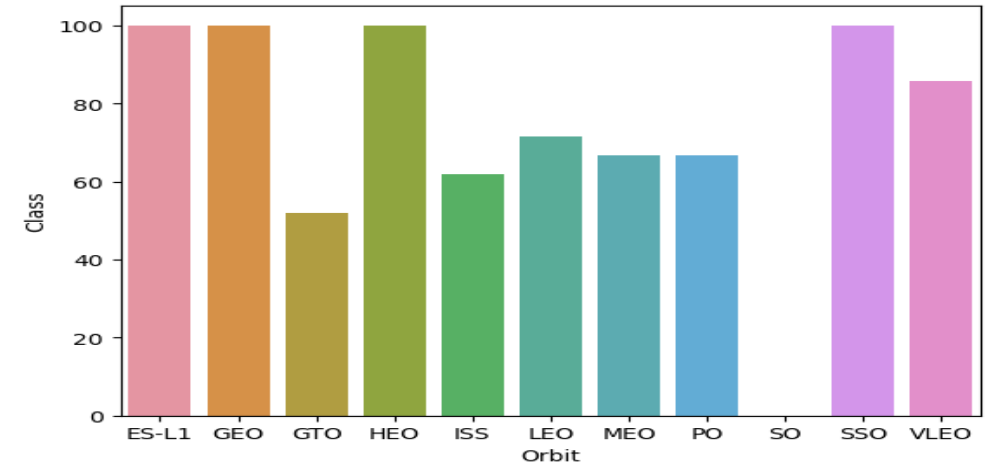
## TASK 2: Visualize the relationship between Payload and Launch Site

We also want to observe if there is any relationship between launch sites and their payload mass.

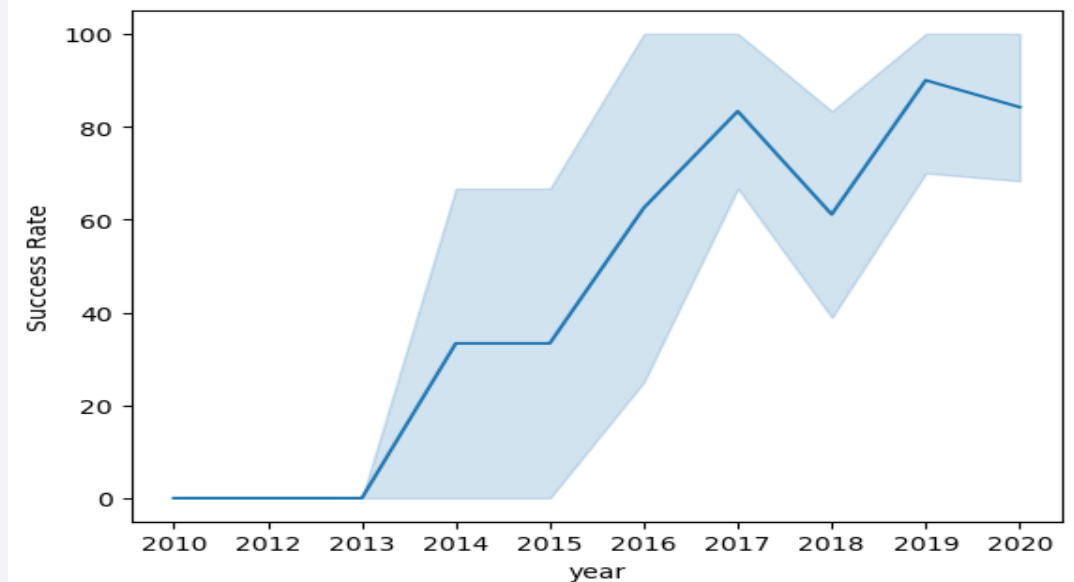
```
# Plot a scatter point chart with x axis to be Pay Load Mass (kg) and y axis to be the launch site, and hue to be the class value
sns.catplot(y="PayloadMass", x="LaunchSite", hue="Class", data=df, aspect = 5)
plt.xlabel("Launch Site", fontsize=20)
plt.ylabel("Pay load Mass (kg)", fontsize=20)
plt.show()
```



Now if you observe Payload Vs. Launch Site scatter point chart you will find for the VAFB-SLC launchsite there are no rockets launched for heavy payload mass (greater than 10000).



Analyze the plotted bar chart try to find which orbits have high success rate.



# EDA with SQL

12

- ▶ We loaded the SpaceX dataset into a PostgreSQL database without leaving the jupyter notebook
- ▶ 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.
- ▶ GitHub URL
- ▶ [https://github.com/Venkatsdandsena/testrepo/blob/a02b100cea689267f5a055115f70292a6b51cf30/EDA%20SQL%20jupyter-labs-eda-sql-coursera\\_sqlite%20\(1\).ipynb](https://github.com/Venkatsdandsena/testrepo/blob/a02b100cea689267f5a055115f70292a6b51cf30/EDA%20SQL%20jupyter-labs-eda-sql-coursera_sqlite%20(1).ipynb)





# Build an Interactive Map with Folium

13

- 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.
- We assigned the feature launch outcomes (failure or success) to class 0 and 1.i.e., 0 for failure, and 1 for success.
- Using the color-labeled marker clusters, we identified which launch sites have relatively high success rate.
- We calculated the distances between a launch site to its proximities. We answered some question for instance: - Are launch sites near railways, highways and coastlines.
- Do launch sites keep certain distance away from cities.
- <https://github.com/Venkatdandsena/testrepo/blob/a02b100cea689267f5a055115f70292a6b51cf30/nteractive%20Visual%20Analytics%20with%20Folium%20lab.ipynb>

# Build a Dashboard with Plotly Dash

---

14

- ▶ 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
- ▶ Add the GitHub URL
- ▶ [https://github.com/Venkatsdandsena/testrepo/blob/a02b100cea689267f5a055115f70292a6b51cf30/spacex\\_dash\\_app%20\(1\).py](https://github.com/Venkatsdandsena/testrepo/blob/a02b100cea689267f5a055115f70292a6b51cf30/spacex_dash_app%20(1).py)

# Predictive Analysis (Classification)

15

- ▶ 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 modelSummarize how you built, evaluated, improved, and found the best performing classification model
- ▶ GitHub
- ▶ <https://github.com/Venkatsena/testrepo/blob/c12c91ea0d76807997a6c79dd32ff153a35233a3/Machine%20Learning%20Prediction.ipynb>

# Results

16

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







Section 2

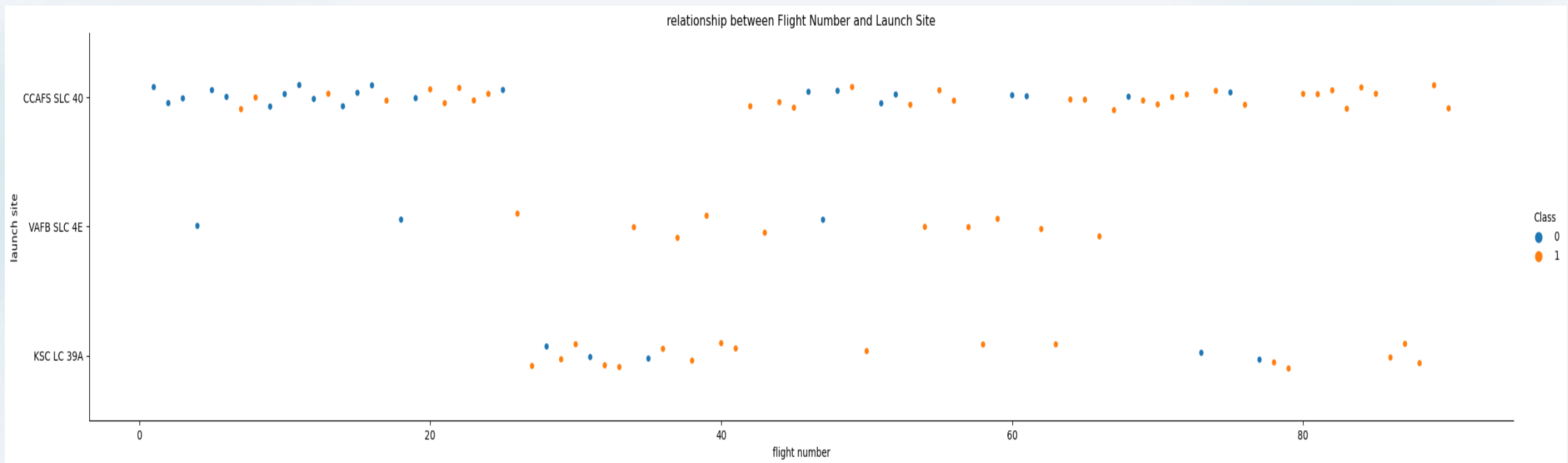
# Insights drawn from EDA



# Flight Number vs. Launch Site

18

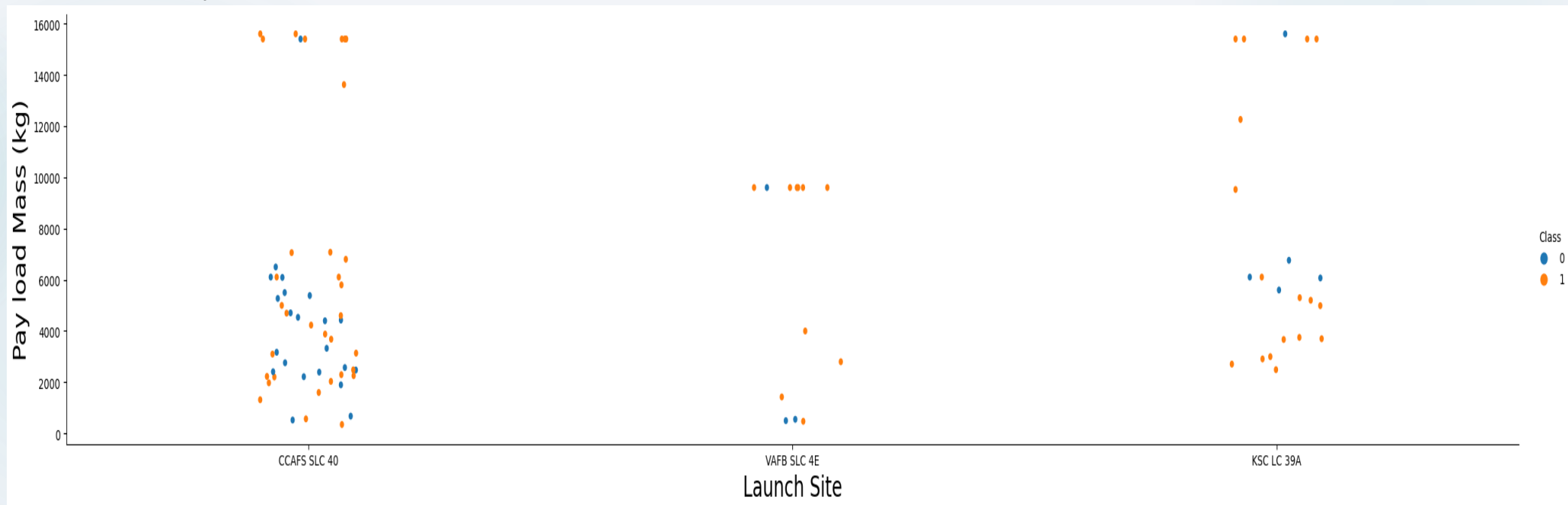
- From the plot, we found that the larger the flight amount at a launch site, the greater the success rate at a launch site.



# Payload vs. Launch Site

19

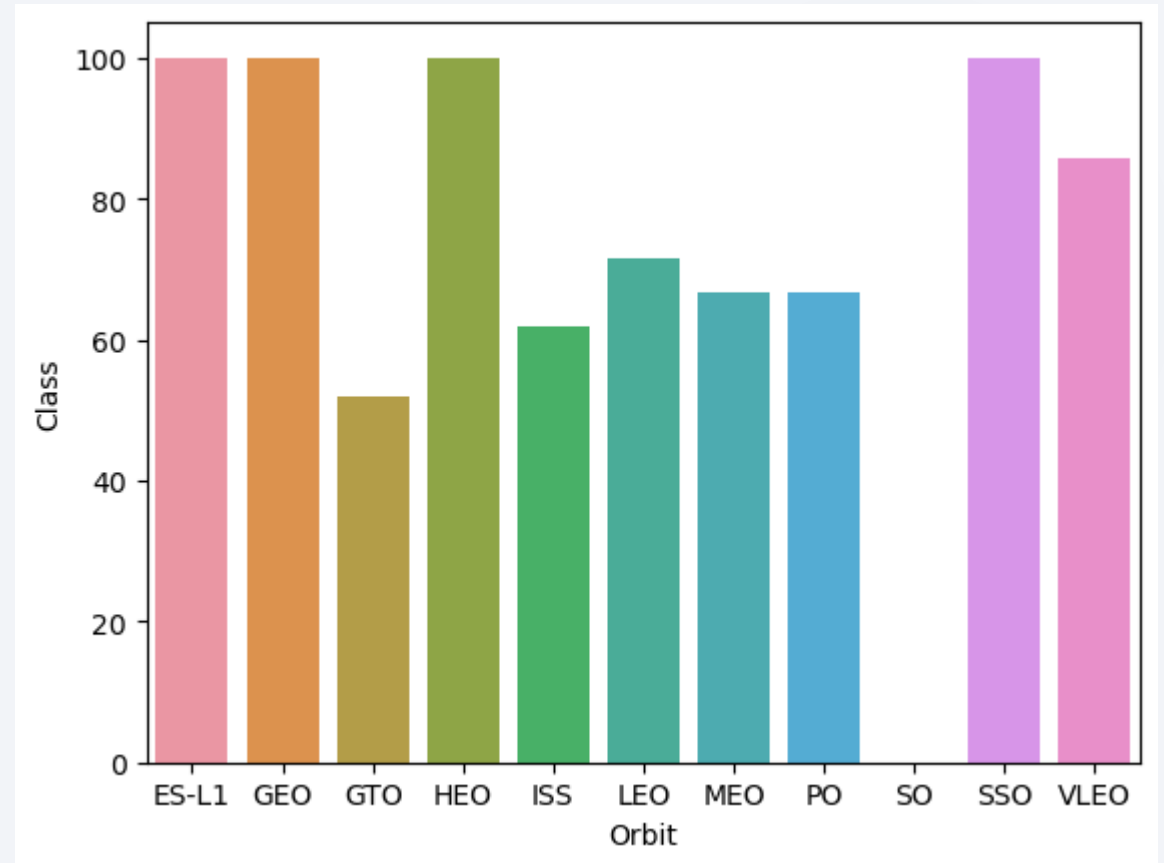
- Payload Vs. Launch Site scatter point chart you will find for the VAFB-SLC launchsite there are no rockets launched for heavypayload mass(greater than 10000).



# Success Rate vs. Orbit Type

20

- ▶ bar chart for the success rate of each orbit type
- ▶ From the plot, we can see that ES-L1, GEO, HEO, SSO, VLEO had the most success rate.

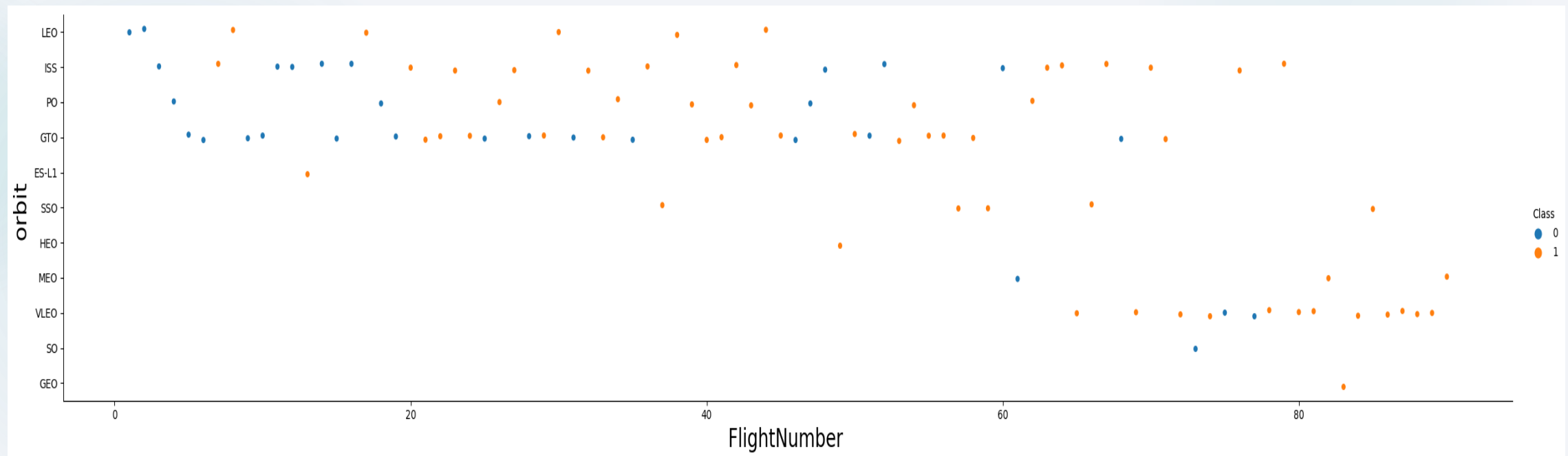




# Flight Number vs. Orbit Type

21

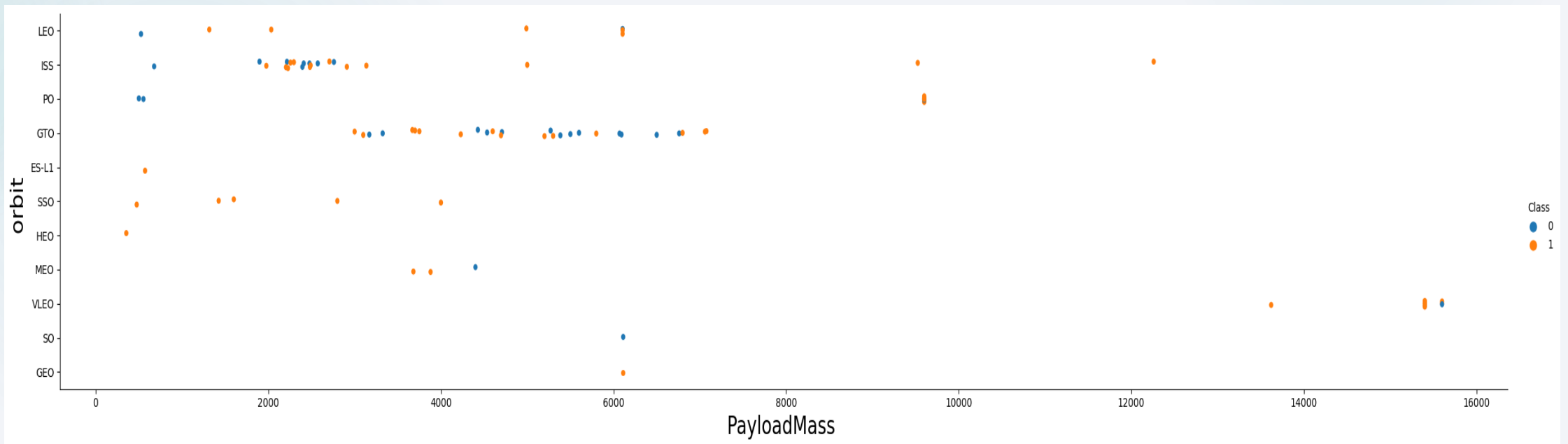
- ▶ scatter point of Flight number vs. Orbit type
- ▶ You should see that in the LEO orbit the Success appears related to the number of flights; on the other hand, there seems to be no relationship between flight number when in GTO orbit.



# Payload vs. Orbit Type

22

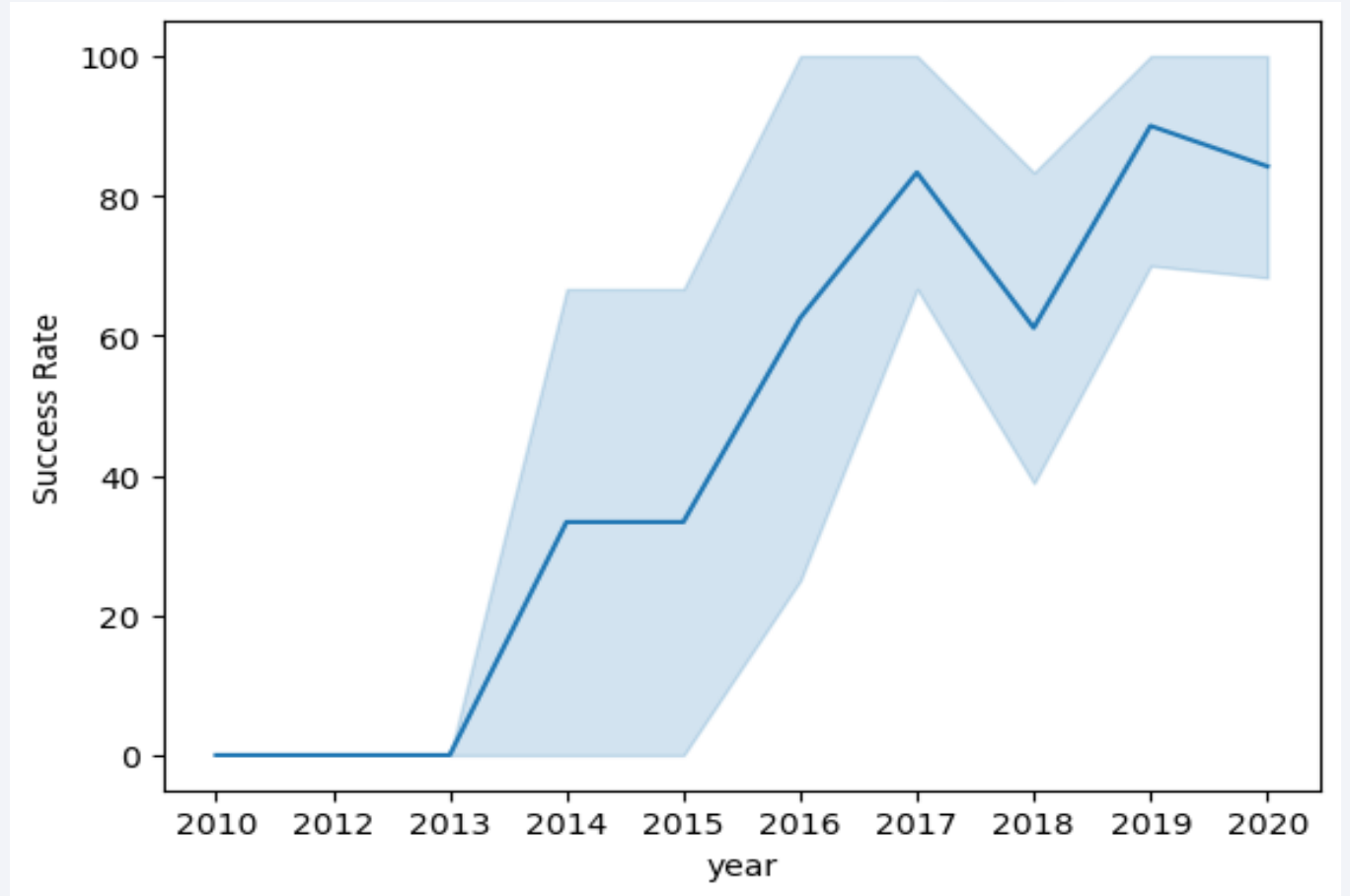
- ▶ scatter point of payload vs. orbit type
- ▶ With heavy payloads the successful landing or positive landing rate are more for Polar, LEO and ISS.
- ▶ However for GTO we cannot distinguish this well as both positive landing rate and negative landing (unsuccessful mission) are both there here.



# Launch Success Yearly Trend

23

- ▶ Show a line chart of yearly average success rate
- ▶ observe that the success rate since 2013 kept increasing till 2020



# All Launch Site Names

24

- ▶ names of the unique launch sites
- ▶ We used the key word DISTINCT to show only unique launch sites from the SpaceX data.

```
%sql select DISTINCT LAUNCH_SITE from SPACEXTbl
```

```
* sqlite:///my_data1.db  
Done.
```

Launch_Site
-------------

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

VAFB SLC-4E
-------------

KSC LC-39A
------------

CCAFS SLC-40
--------------



# Launch Site Names Begin with 'CCA'

25

## ► 5 records where launch sites begin with 'CCA'

### Task 2

Display 5 records where launch sites begin with the string 'CCA'

```
%sql select * from SPACEXTBL where launch_site like 'CCA%' limit 5
```

```
* sqlite:///my_data1.db
```

Done.

Date	Time (UTC)	Booster_Version	Launch_Site	Payload	PAYLOAD_MASS_KG_	Orbit	Customer	Mission_Outcome	Landing_Outcome
04-06-2010	18:45:00	F9 v1.0 B0003	CCAFS LC-40	Dragon Spacecraft Qualification Unit	0	LEO	SpaceX	Success	Failure (parachute)
08-12-2010	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)
22-05-2012	07:44:00	F9 v1.0 B0005	CCAFS LC-40	Dragon demo flight C2	525	LEO (ISS)	NASA (COTS)	Success	No attempt
08-10-2012	00:35:00	F9 v1.0 B0006	CCAFS LC-40	SpaceX CRS-1	500	LEO (ISS)	NASA (CRS)	Success	No attempt
01-03-2013	15:10:00	F9 v1.0 B0007	CCAFS LC-40	SpaceX CRS-2	677	LEO (ISS)	NASA (CRS)	Success	No attempt

# Total Payload Mass

26

- ▶ the total payload carried by boosters from NASA
- ▶ We calculated the total payload carried by boosters from NASA as 45596 using the query below.

## Task 3

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

```
%sql SELECT SUM(      payload_mass__kg_) from spacextbl where Customer = 'NASA (CRS)'
```

```
* sqlite:///my_data1.db  
Done.
```

SUM( payload_mass__kg_)
45596

# Average Payload Mass by F9 v1.1

27

- ▶ average payload mass carried by booster version F9 v1.1
- ▶ We calculated the average payload mass carried by booster version F9 v1.1 as 2928.4

## Task 4

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

```
[43]: %sql select AVG(payload_mass__kg_) as average FROM SPACEXTBL WHERE booster_version like 'F9 v1.1%'
```

```
* sqlite:///my_data1.db  
Done.
```

```
[43]: 

| average            |
|--------------------|
| 2534.6666666666665 |


```

# First Successful Ground Landing Date

28

- ▶ dates of the first successful landing outcome on ground pad

## Task 5

List the date when the first succesful landing outcome in ground pad was acheived.

*Hint: Use min function*

```
5]: %sql select min(date) as Date from SPACEXTbl where mission_outcome like 'Success'
```

```
* sqlite:///my_data1.db  
Done.
```

```
5]: 

| Date       |
|------------|
| 01-03-2013 |


```

## Successful Drone Ship Landing with Payload between 4000 and 6000 29

- List the names of boosters which have successfully landed on drone ship and had payload mass greater than 4000 but less than 6000

### Task 6

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

```
In [63]: %sql select booster_version from SPACEXTbl where (mission_outcome like 'Success') AND (PAYLOAD_MASS_KG_ BETWEEN 4000 AND 6000) AND ("landing_outcome
```

```
* sqlite:///my_data1.db  
Done.
```

```
Out[63]: Booster_Version
```

```
F9 FT B1022
```

```
F9 FT B1026
```

```
F9 FT B1021.2
```

```
F9 FT B1031.2
```



# Total Number of Successful and Failure Mission Outcomes

30

- ▶ the total number of successful and failure mission outcomes

## Task 7

List the total number of successful and failure mission outcomes

```
%sql SELECT mission_outcome, count(*) as Count FROM SPACEXTbl GROUP by mission_outcome ORDER BY mission_outcome
```

```
* sqlite:///my_data1.db
```

```
Done.
```

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

# Boosters Carried Maximum Payload

31

- List the names of the booster which have carried the maximum payload mass

## Task 8

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

```
]]: maxm = %sql select max(payload_mass__kg_) from SPACEXTBL
maxv = maxm[0][0]
%sql select booster_version from SPACEXTBL where payload_mass__kg_=(select max(payload_mass__kg_) from SPACEXTBL)

* sqlite:///my_data1.db
Done.
* sqlite:///my_data1.db
Done.

]: 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

32

- ▶ List the failed landing\_outcomes in drone ship, their booster versions, and launch site names for in year 2015
- ▶ We used a combinations of the WHERE clause, LIKE, AND, and BETWEEN conditions to filter for failed landing outcomes in drone ship, their booster versions, and launch site names for year 2015

## Task 9

List the records which will display the month names, failure landing\_outcomes in drone ship ,booster versions, launch\_site for the months in year 2015.

**Note: SQLite does not support monthnames. So you need to use substr(Date, 4, 2) as month to get the months and substr(Date,7,4)='2015' for year.**

```
5]: %sql select SUBSTR(DATE,4,2) as Month, "landing__outcome", booster_version, launch_site from SPACEXTBL where DATE like '%2015' AND ("landing__outcome"
```

```
* sqlite:///my_data1.db  
Done.
```

```
5]:
```

	Month	"landing__outcome"	Booster_Version	Launch_Site
	01	landing__outcome	F9 v1.1 B1012	CCAFS LC-40
	04	landing__outcome	F9 v1.1 B1015	CCAFS LC-40

# 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

## Task 10

Rank the count of successful landing\_outcomes between the date 04-06-2010 and 20-03-2017 in descending order.

```
)]: %sql select "landing_outcome", count(*) as count from SPACEXTBL where Date >= '04-06-2010' AND Date <= '20-03-2017' GROUP by "landing_outcome" ORDER  
* sqlite:///my_data1.db  
Done.
```

```
)]:
```

Landing_Outcome	count
Success	20
No attempt	10
Success (drone ship)	8
Success (ground pad)	6
Failure (drone ship)	4
Failure	3
Controlled (ocean)	3
Failure (parachute)	2
No attempt	1



Section 3

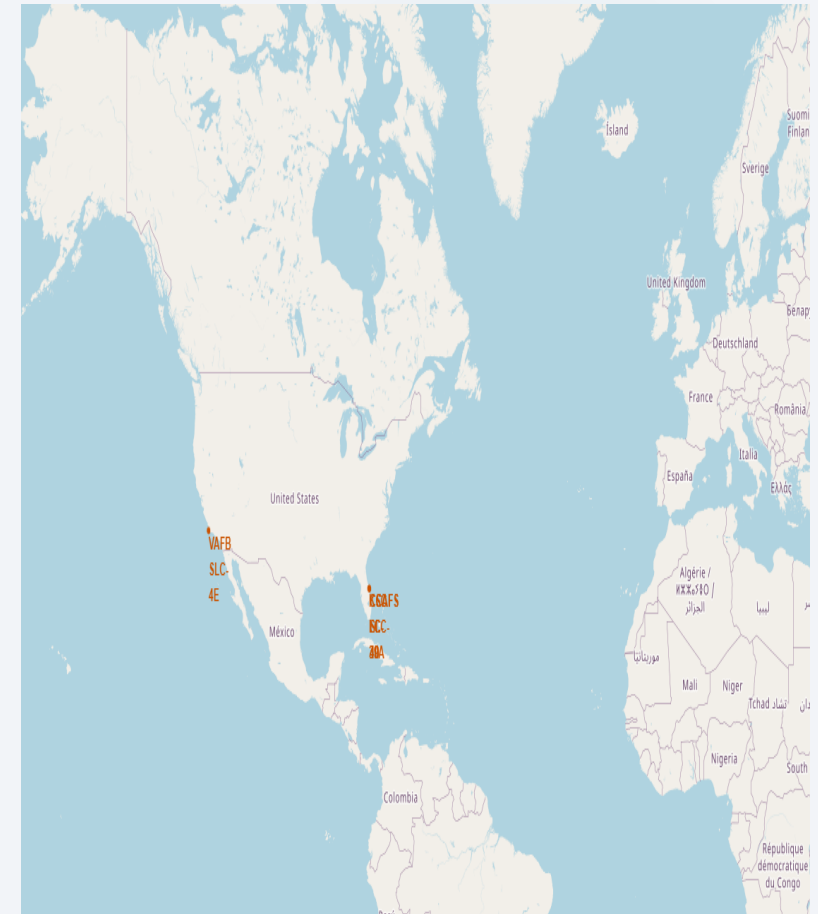
# Launch Sites Proximities Analysis



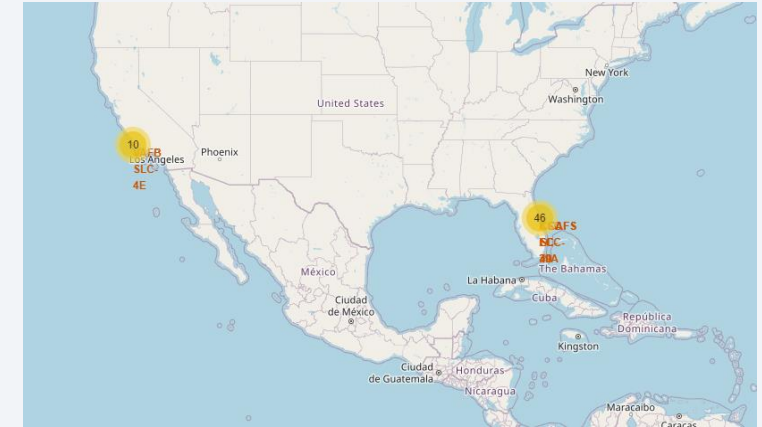
# All the launch sites

35

Findings from the map is that all the launch site are in United state And near to the sea or water.

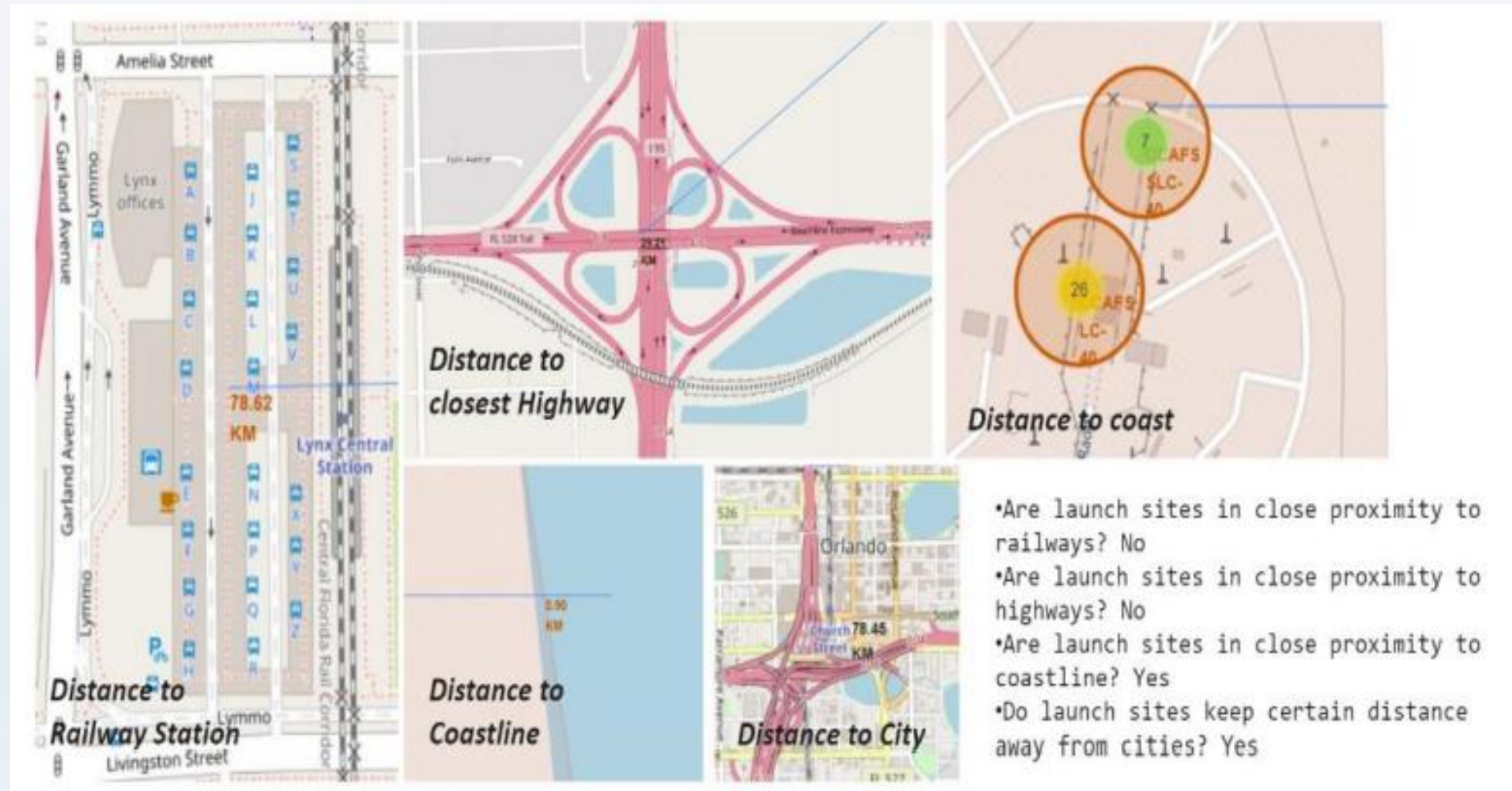


- ▶ Marker clusters can be a good way to simplify a map containing many markers having the same coordinate.
- ▶ From the color-labeled markers in marker clusters, easily identify which launch sites have relatively high success rates



# Launch site distance from railways, coast and city

37







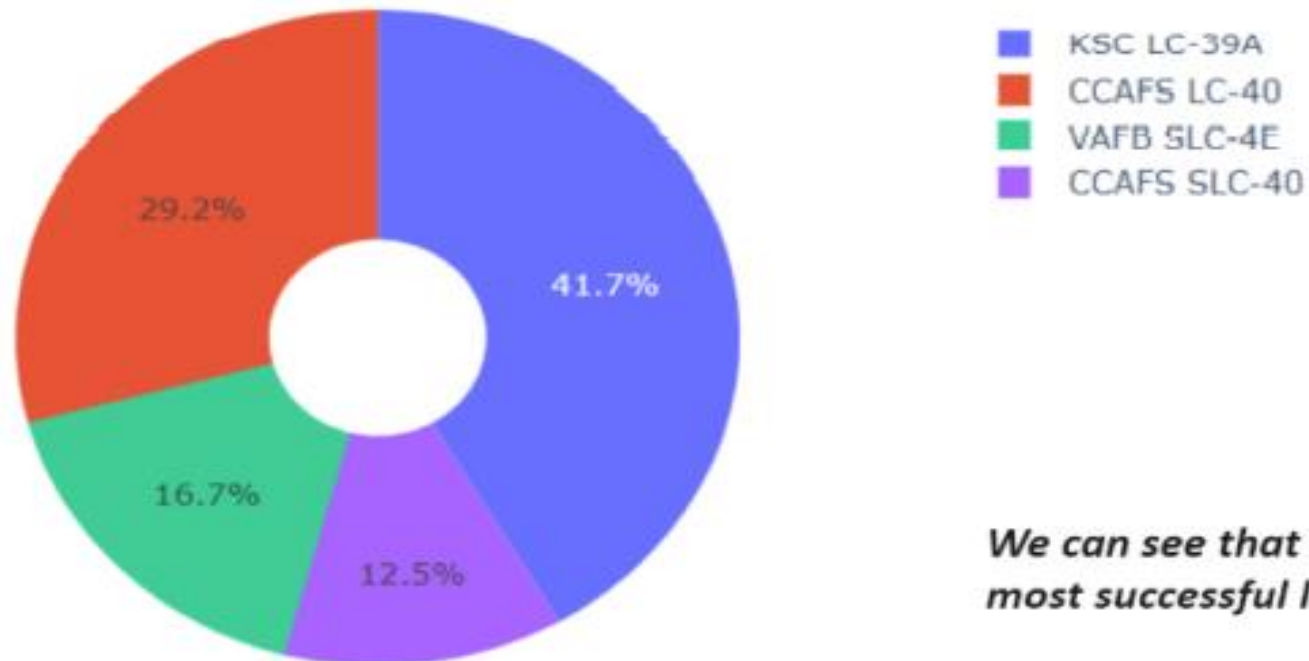
Section 4

# Build a Dashboard with Plotly Dash

# Total launch from all sites

39

Total Success Launches By all sites

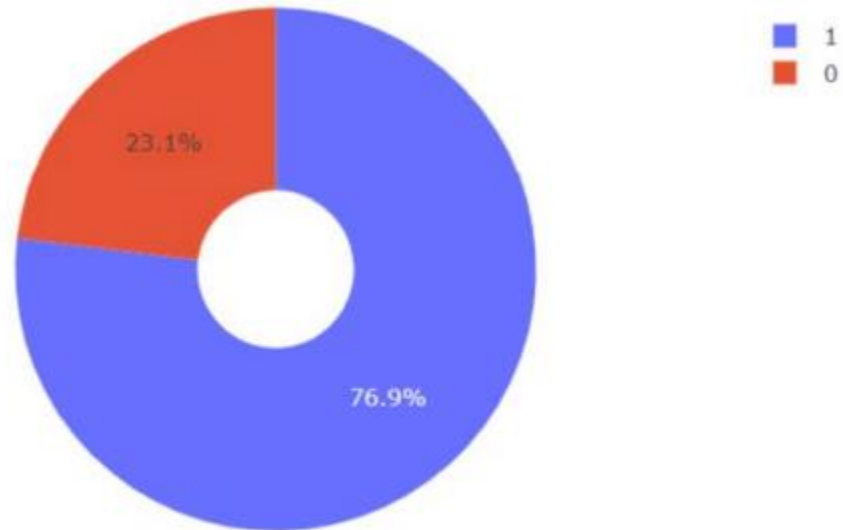


*We can see that KSC LC-39A had the most successful launches from all the sites*



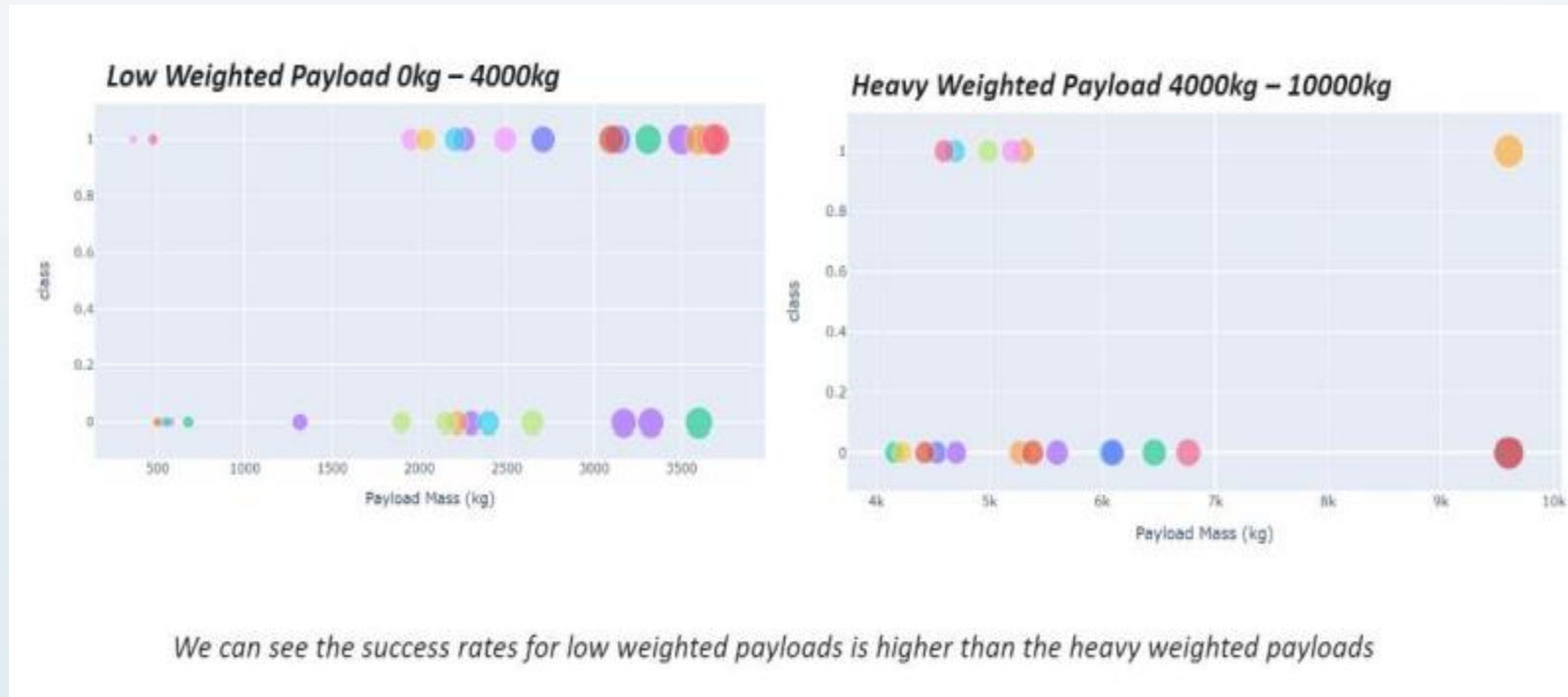
## Pie chart showing the Launch site with the highest launch success ratio

40



*KSC LC-39A achieved a 76.9% success rate while getting a 23.1% failure rate*

## Scatter plot of Payload vs Launch Outcome for all sites, with different payload selected in the range slide



The background of the slide is an abstract composition. On the left, there is a solid blue area. To the right, a series of white and light blue curved lines create a sense of motion and depth, resembling a tunnel or a futuristic architectural space. A bright red rectangle is positioned in the upper right corner.

Section 5

# Predictive Analysis (Classification)

- The decision tree classifier is the model with the highest classification accuracy

Find the method performs best:

```
: 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 0.8732142857142856

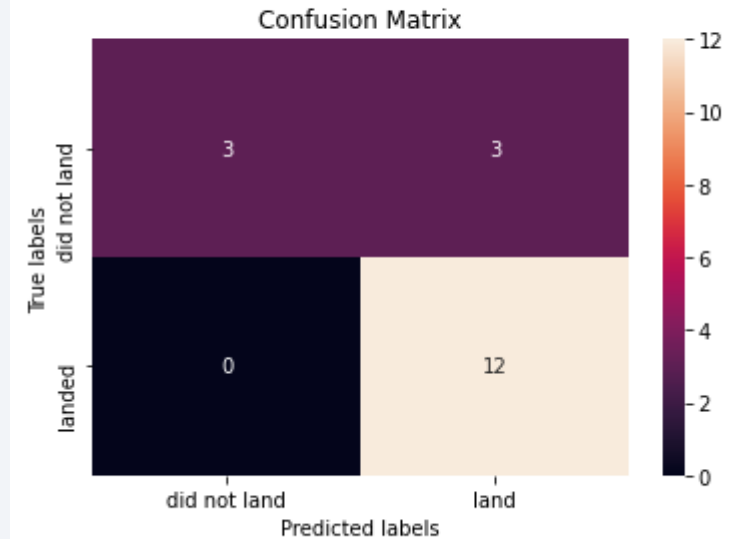
Best params is : {'criterion': 'gini', 'max\_depth': 6, 'max\_features': 'auto', 'min\_samples\_leaf': 2, 'min\_samples\_split': 5, 'splitter': 'random'}

# Confusion Matrix

44

The confusion matrix for the decision tree classifier shows that the classifier can distinguish between the different classes. The major problem is the false positives .i.e., unsuccessful landing marked as successful landing by the classifier.

```
yhat_tree = tree_cv.predict(X_test)
plot_confusion_matrix(Y_test, yhat_tree)
```





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

Thank you!

