



IBM Developer  
SKILLS NETWORK

# Winning Space Race with Data Science

Karan Verma  
19<sup>th</sup> March, 2022



# Outline

---

- ❖ 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

- 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

- 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 needs to be in place to ensure a successful landing program.



Section 1

# Methodology

# 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.

❖ For reference:

<https://github.com/karannxo/Caps-tone-Project-SPACE-Y/blob/f12c6ce17fda751b612dcae515775e2b9443e055/Week%201%20API.ipynb>

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

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

Check the content of the response

In [8]: print(response.content)

In [9]: static_json_url="https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBM-DS0321EN-SkillsNetwork/datasets/API_call_spacex_api.json"

We should see that the request was successful with the 200 status response code

In [10]: response.status_code

Out[10]: 200

Now we decode the response content as a json using .json() and turn it into a Pandas dataframe using .json_normalize()

In [11]: # Use json_normalize method to convert the json result into a dataframe
data = pd.json_normalize(response.json())

Using the dataframe 'data' print the first 5 rows

In [12]: # Get the head of the dataframe
data.head()

In [27]: # Calculate the mean value of PayloadMass column
pm_mean = data_falcon9['PayloadMass'].mean()
# Replace the np.nan values with its mean value
data_falcon9['PayloadMass'] = data_falcon9['PayloadMass'].replace(np.nan, pm_mean)
data_falcon9
```



# Data Collection - Scrapping

- ❖ We applied web scrapping to webscrap Falcon 9 launch records with BeautifulSoup
- ❖ We parsed the table and converted it into a pandas dataframe.
- ❖ For reference:  
<https://github.com/karannxo/Caps-tone-Project-SPACE-Y/blob/f12c6ce17fda751b612dcae515775e2b9443e055/Week%201%20Web.ipynb>

```
In [4]: 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.

In [5]: # use requests.get() method with the provided static_url
# assign the response to a object
data = requests.get(static_url)

Create a BeautifulSoup object from the HTML response

In [6]: # Use BeautifulSoup() to create a BeautifulSoup object from a response text content
soup = BeautifulSoup(data.text, 'html.parser')

<td></td></tr>
</th></tr>

Next, we just need to iterate through the <th> elements and apply the provided extract_column_from_header() to extract column name one by one

In [10]: column_names = []
for i in first_launch_table.find_all('th'):
    if extract_column_from_header(i) != None:
        if len(extract_column_from_header(i)) > 0:
            column_names.append(extract_column_from_header(i))

# Apply find_all() function with 'th' element on first_launch_table
# Iterate each th element and apply the provided extract_column_from_header() to get a column name
# Append the Non-empty column name ('if name is not None and len(name) > 0') into a list called column_names

Check the extracted column names

In [11]: print(column_names)
```

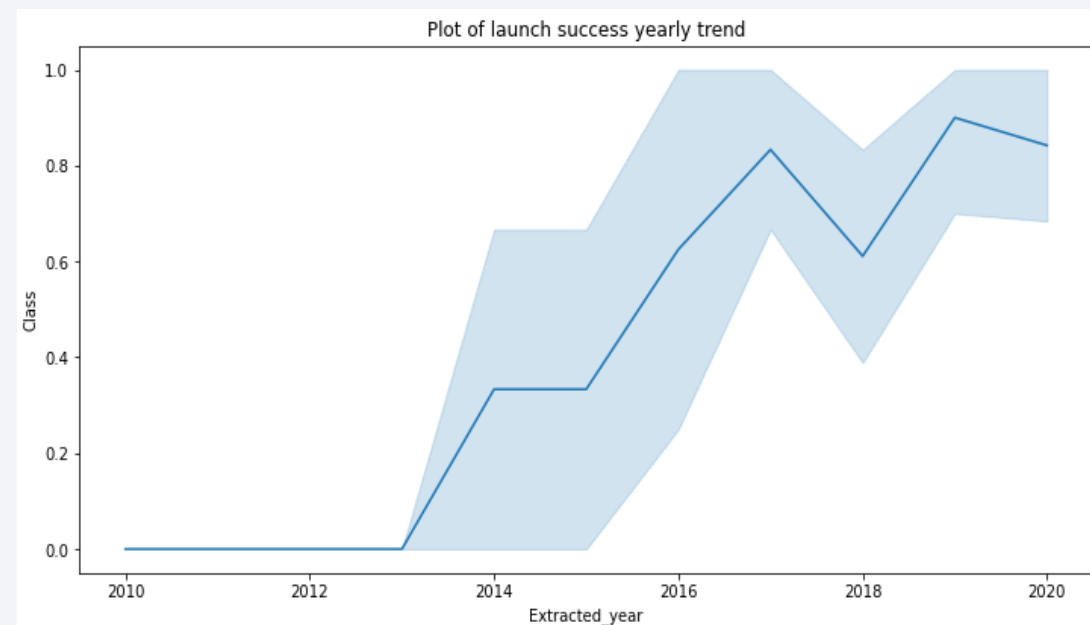
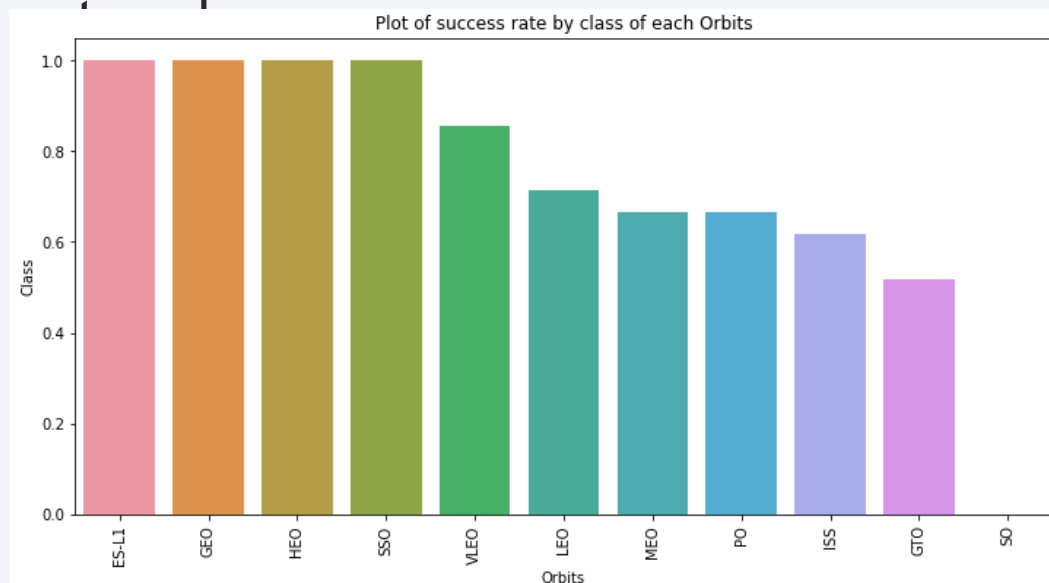
# Data Wrangling



- ❖ 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.
- ❖ For reference:  
<https://github.com/karannxo/Capstone-Project-SPACE-Y/blob/f12c6ce17fda751b612dcae515775e2b9443e055/Week%201%20Data%20Wrangling.ipynb>

# 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



- For reference:  
<https://github.com/karannxo/Capstone-Project-SPACE-Y/blob/f12c6ce17fda751b612dcae515775e2b9443e055/Week%202%20EDA%20Python.ipynb>

# EDA with SQL

---

- ❖ 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.
- ❖ For reference: <https://github.com/karannxo/Capstone-Project-SPACE-Y/blob/f12c6ce17fda751b612dcae515775e2b9443e055/Week%20%20EDA%20SQL.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.
- ❖ 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.
- ❖ For reference : <https://github.com/karannxo/Capstone-Project-SPACE-Y/blob/f12c6ce17fda751b612dcae515775e2b9443e055/Week%203%20Dashboard.ipynb>



# 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.
- ❖ For reference: [https://github.com/karannxo/Capstone-Project-SPACE-Y/blob/9322739d34ad31ea2ec124f15154d18001dc149a/spacex\\_dash\\_app.py](https://github.com/karannxo/Capstone-Project-SPACE-Y/blob/9322739d34ad31ea2ec124f15154d18001dc149a/spacex_dash_app.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.
- ❖ For reference: <https://github.com/karannxo/Capstone-Project-SPACE-Y/blob/9322739d34ad31ea2ec124f15154d18001dc149a/Week%204%20ML.ipynb>

# Results

---

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



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 are layered over a faint, dark grid pattern, creating a sense of depth and movement.

Section 2

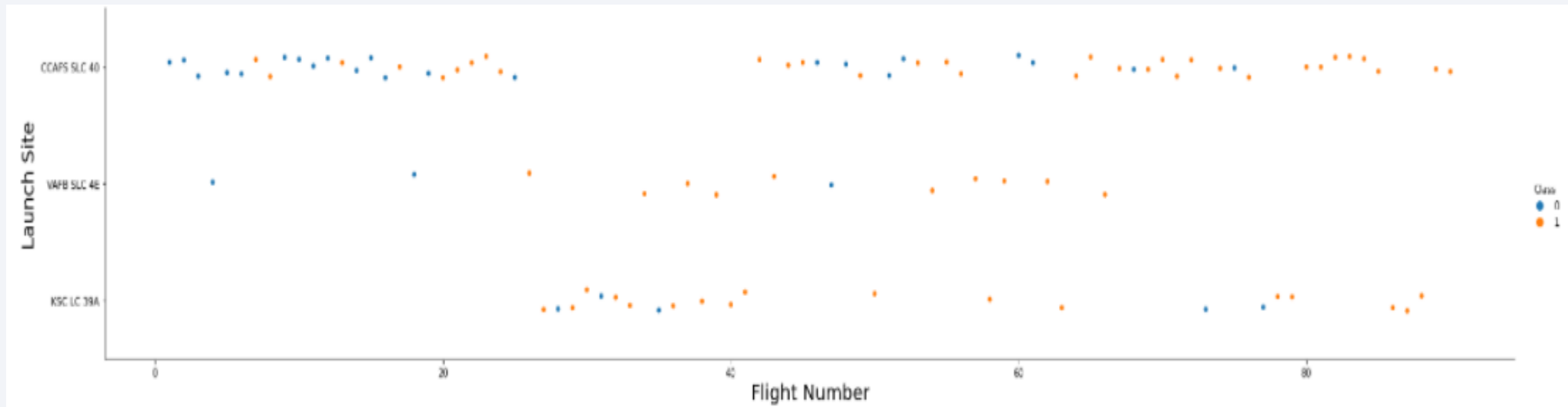
# Insights drawn from EDA



# 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.





## Payload vs. Launch Site



The greater the payload mass for launch site CCAFS SLC 40 the higher the success rate for the rocket.



# Success Rate vs. Orbit Type

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

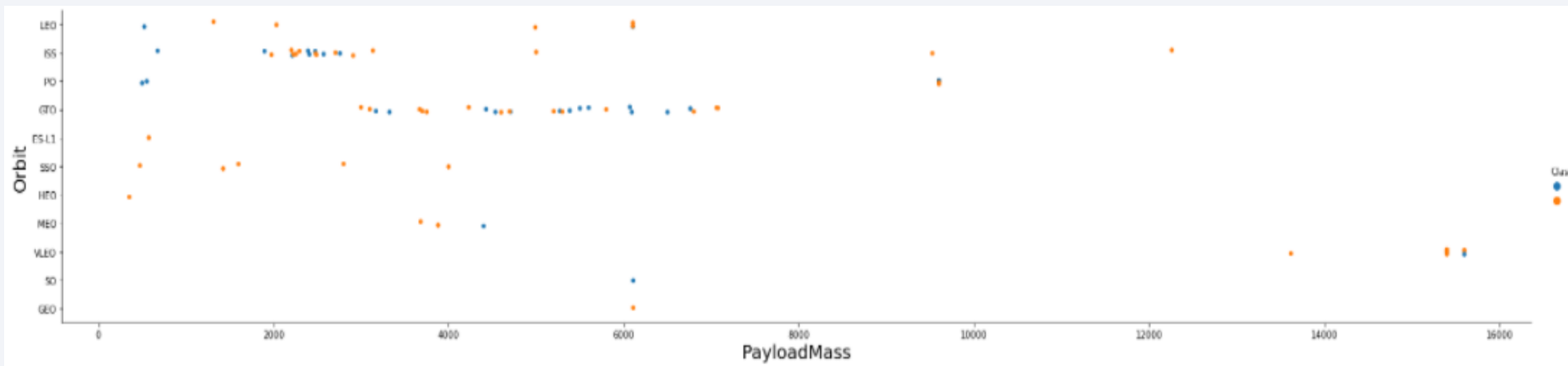




# Payload vs. Orbit Type

---

- ❖ We can observe that with heavy payloads, the successful landing are more for PO, LEO and ISS orbits.



# Launch Success Yearly Trend

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





# All Launch Site Names

❖ We used the key word **DISTINCT** to show only unique launch sites from the SpaceX data.

## Task 1

Display the names of the unique launch sites in the space mission

In [9]:

```
%%sql  
select distinct(launch_site) as "Diff launch" from SPACEXTBL;
```

```
* ibm_db_sa://zxj01478:***@824dfd4d-99de-440d-9991-629c01b3832d.bs2io90108kqb1od8lcg.databases.appdomain.cloud:30119/bludb  
Done.
```

Out[9]: Diff launch

CCAFS LC-40

CCAFS SLC-40

KSC LC-39A

VAFB SLC-4E

# Launch Site Names Begin with 'CCA'

## Task 2

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

In [10]:

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

```
* ibm_db_sa://zxj01478:***@824dfd4d-99de-440d-9991-629c01b3832d.bs2io90l08kqb1od8lcg.databases.appdomain.cloud:30119/blddb
Done.
```

Out[10]:

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

❖ We used the query above to display 5 records where launch sites begin with 'CCA'

# Total Payload Mass

---

- ❖ 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)

```
In [11]: %sql select sum(payload_mass__kg_) as "Total" from SPACEXTBL where customer = 'NASA (CRS)'
```

```
* ibm_db_sa://zxj01478:***@824dfd4d-99de-440d-9991-629c01b3832d.bs2io90l08kqb1od8lcg.databases.appdomain.cloud:30119/bludb  
Done.
```

```
Out[11]: Total
```

```
45596
```

# Average Payload Mass by F9 v1.1

---

- ❖ We calculated the average payload mass carried by booster version F9 v1.1 as 2534 KG.

## Task 4

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

In [13]:

```
%%sql
select avg(payload_mass__kg_) as AVG from SPACEXTBL where booster_version like 'F9 v1.1%'
```

```
* ibm_db_sa://zxj01478:***@824dfd4d-99de-440d-9991-629c01b3832d.bs2io90l08kqb1od8lcg.databases.appdomain.cloud:30119/bludb
Done.
```

Out[13]:

**AVG**

2534

# First Successful Ground Landing Date

---

- ❖ We observed that the dates of the first successful landing outcome on ground pad was 22<sup>nd</sup> December 2015.

## Task 5

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

*Hint: Use min function*

```
In [24]: %sql select min(date) as Date from SPACEXTBL where "Landing _Outcome" like 'Success (ground pad)'
```

\* ibm\_db\_sa://zxj01478:\*\*\*@824dfd4d-99de-440d-9991-629c01b3832d.bs2io90l08kqb1od8lcg.databases.appdomain.cloud:30119/bludb  
Done.

```
Out[24]:      DATE  
2015-12-22
```



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

---

- ❖ We used the **WHERE** clause to filter for boosters which have successfully landed on drone ship and applied the **AND** condition to determine successful landing with 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 [25]:

```
%%sql
select booster_version from SPACEXTBL where ("Landing _Outcome" like 'Success (drone ship)') and (payload_mass__kg_ BETWEEN 4000 AND 6000)
```

```
* ibm_db_sa://zxj01478:***@824dfd4d-99de-440d-9991-629c01b3832d.bs2io90l08kqb1od8lcg.databases.appdomain.cloud:30119/bludb
Done.
```

Out[25]: **booster\_version**

F9 FT B1022

F9 FT B1026

F9 FT B1021.2

F9 FT B1031.2

# Total Number of Successful and Failure Mission Outcomes

---

❖ We used wildcard like '%' to filter for **WHERE** Mission Outcome was a success or a failure.

## Task 7

List the total number of successful and failure mission outcomes

In [26]:

```
%sql SELECT mission_outcome, count(*) as Count FROM SPACEXTBL GROUP by mission_outcome
```

```
* ibm_db_sa://zxj01478:***@824dfd4d-99de-440d-9991-629c01b3832d.bs2io90l08kqb1od8l1cg.databases.appdomain.cloud:30119/bludb  
Done.
```

Out[26]:

mission_outcome	COUNT
Failure (in flight)	1
Success	99
Success (payload status unclear)	1

# Boosters Carried Maximum Payload

- ❖ We determined the booster that have carried the maximum payload using a subquery in the **WHERE** clause and the **MAX()** function.

## Task 8

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

In [28]:

```
%%sql
select booster_version from SPACEXTBL
where payload_mass_kg_=(select max(payload_mass_kg_) from SPACEXTBL)
```

```
* ibm_db_sa://zxj01478:***@824dfd4d-99de-440d-9991-629c01b3832d.bs2io90108kqb1od81cg.databases.appdomain.cloud:30119/bludb
Done.
```

Out[28]:

**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

---

- ❖ We used a combinations of the **WHERE** clause, **LIKE**, and **AND** conditions to filter for failed landing outcomes in drone ship, their booster versions, and launch site names for year 2015.

## Task 9

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

In [34]:

```
%%sql
select ("Landing _Outcome") as Landing_outcome, booster_version, launch_site from SPACEXTBL
where year(DATE) = 2015 AND "Landing _Outcome" like 'Failure (drone ship)'
```

```
* ibm_db_sa://zxj01478:***@824dfd4d-99de-440d-9991-629c01b3832d.bs2io90l08kqb1od8lcg.databases.appdomain.cloud:30119/bludb
Done.
```

Out[34]:

landing_outcome	booster_version	launch_site
Failure (drone ship)	F9 v1.1 B1012	CCAFS LC-40
Failure (drone ship)	F9 v1.1 B1015	CCAFS LC-40

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

- ❖ We selected Landing outcomes and the **COUNT** of landing outcomes from the data and used the **WHERE** clause to filter for landing outcomes **BETWEEN** 2010-06-04 to 2010-03-20.
- ❖ We then applied the **GROUP BY** clause to group the landing outcomes and the **ORDER BY** clause to order the grouped landing outcome in descending order.

### Task 10

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 [37]:

```
%%sql
select ("Landing _Outcome") as Landing_outcome, count(*) as count from SPACEXTBL
where Date BETWEEN '2010-06-04' AND '2017-03-20'
GROUP by "Landing _Outcome" ORDER BY count Desc
```

```
* ibm_db_sa://zxj01478:***@824dfd4d-99de-440d-9991-629c01b3832d.bs2io90l08kqb1od8lcg.databases.appdomain.cloud:30119/bludb
Done.
```

Out[37]:

landing_outcome	COUNT
No attempt	10
Failure (drone ship)	5
Success (drone ship)	5
Controlled (ocean)	3
Success (ground pad)	3
Failure (parachute)	2
Uncontrolled (ocean)	2
Precluded (drone ship)	1

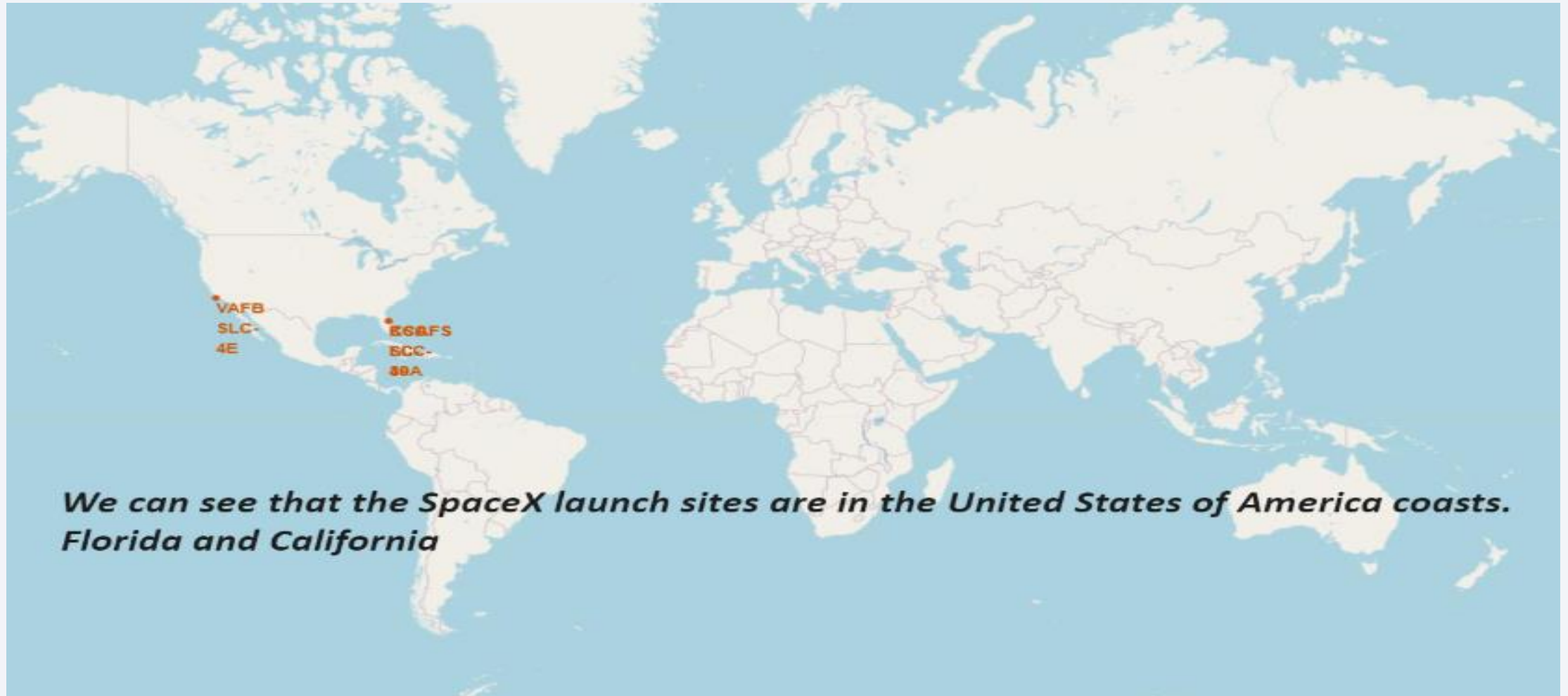
Section 4

# Launch Sites Proximities Analysis



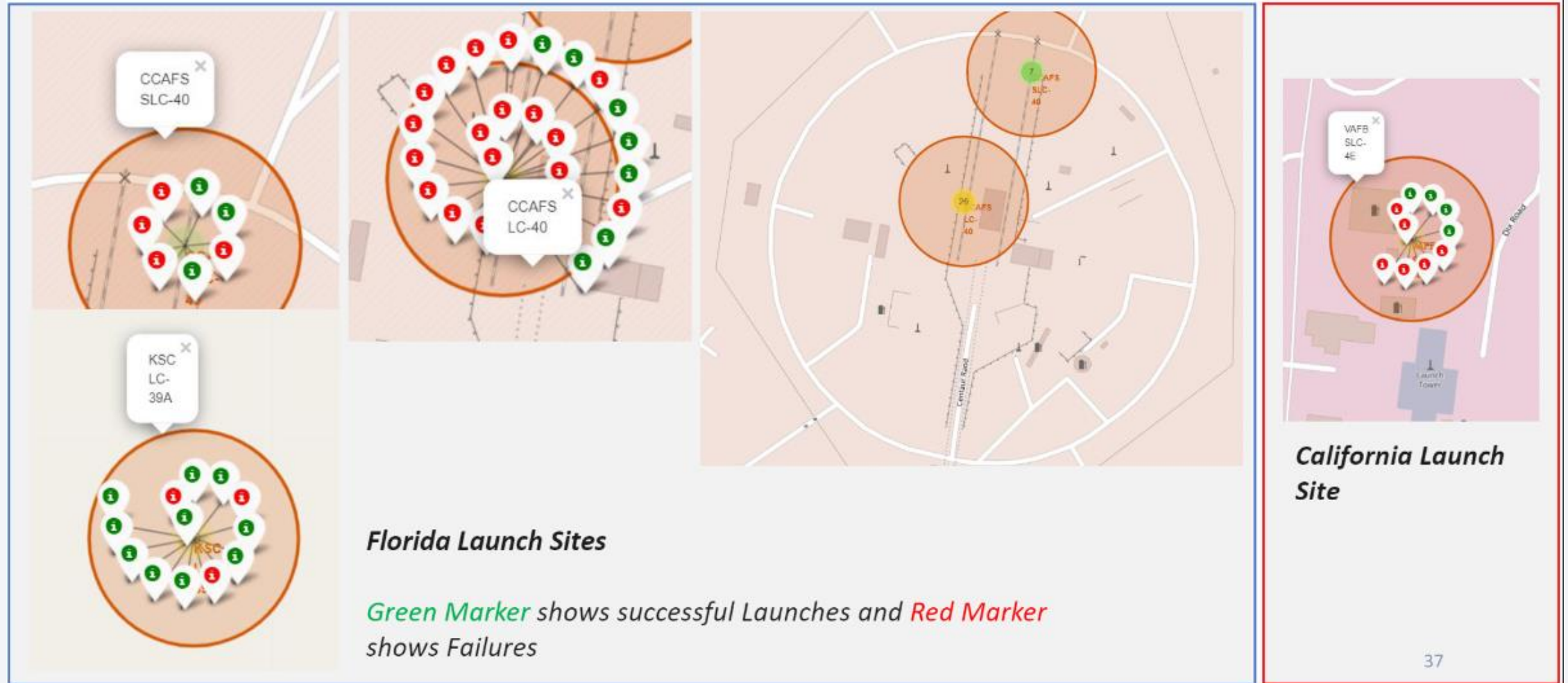
# All launch sites global map markers

---

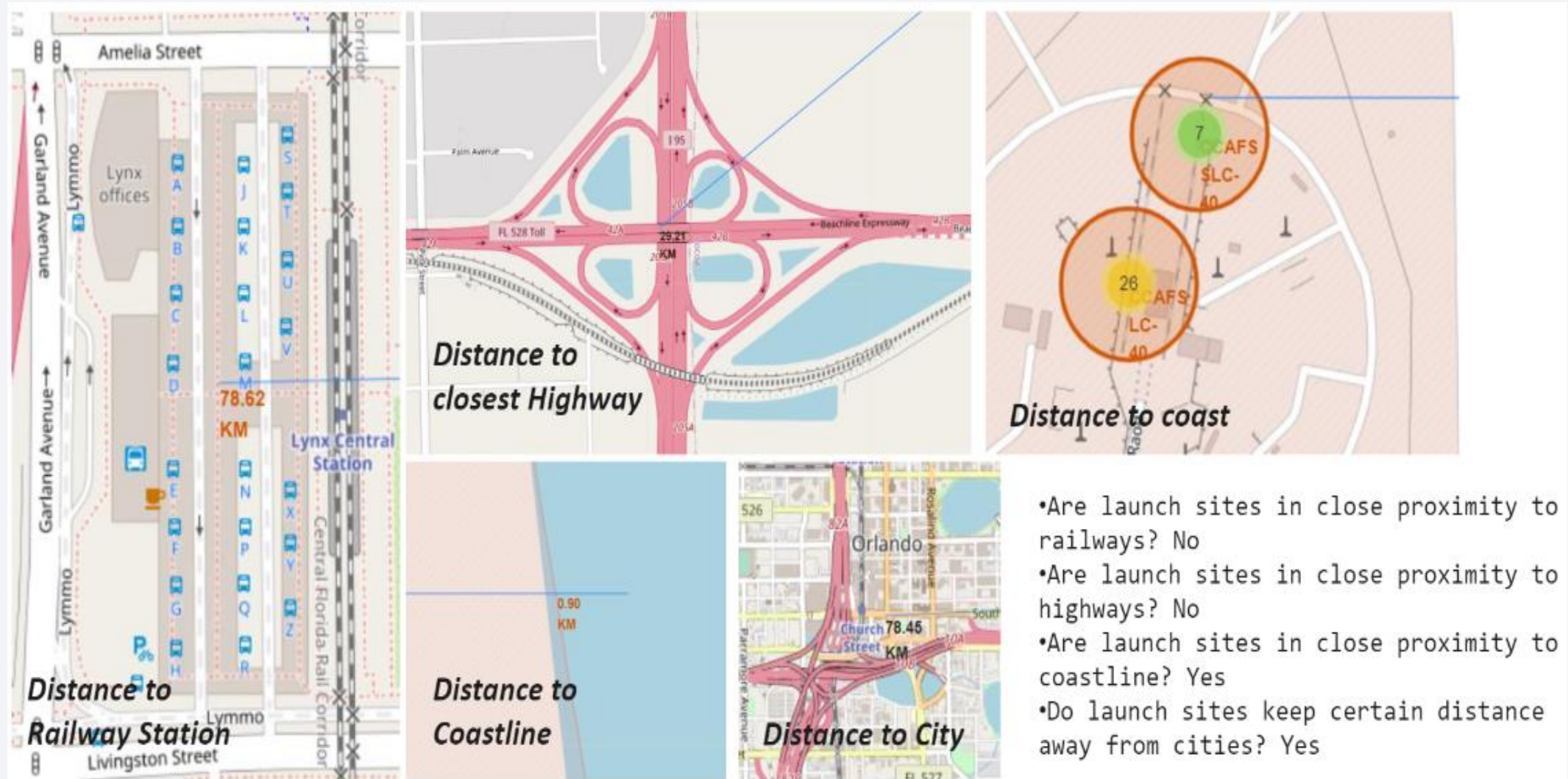




# Markers showing launch sites with color labels



# Launch Site distance to landmarks







Section 5

# Build a Dashboard with Plotly Dash

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

---

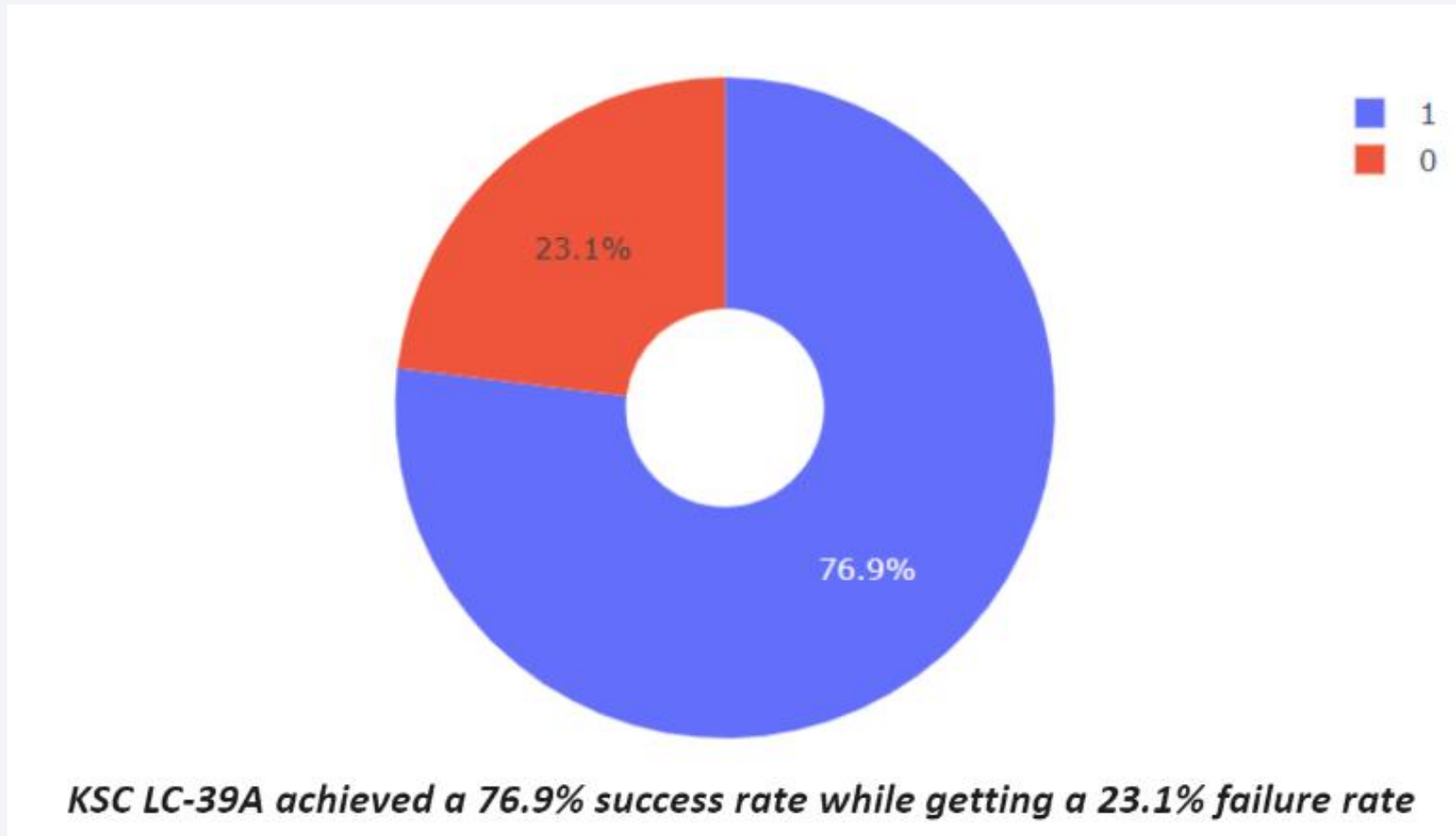
Total Success Launches By all sites



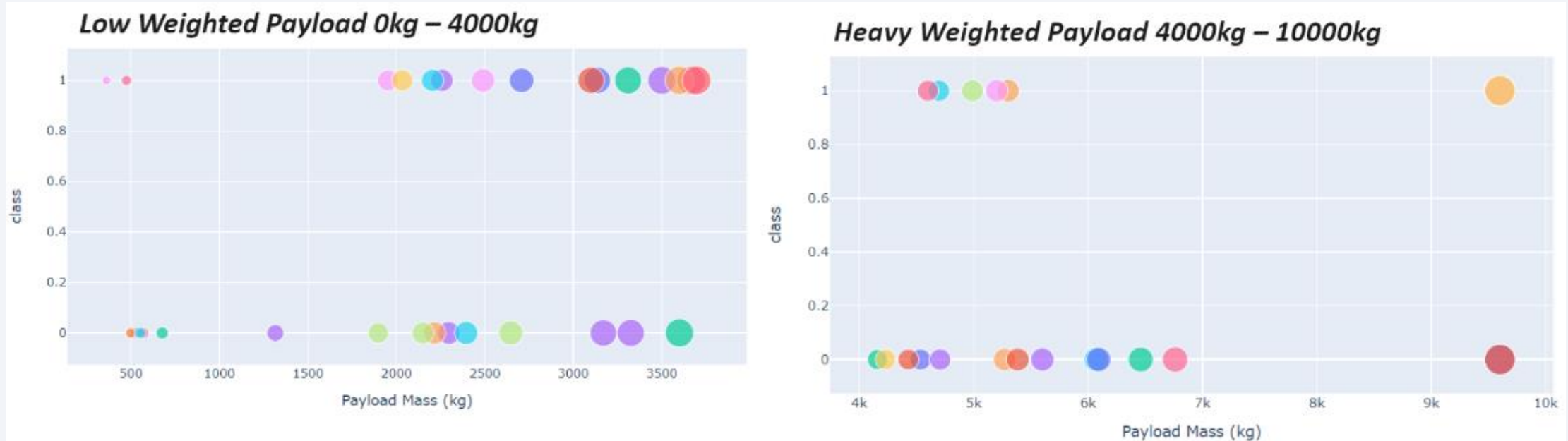
***We can see that KSC LC-39A had the most successful launches from all the sites***

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

---



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



*We can see the success rates for low weighted payloads is higher than the heavy weighted payloads*



Section 6

# Predictive Analysis (Classification)



# Classification Accuracy

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

```
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_)
```

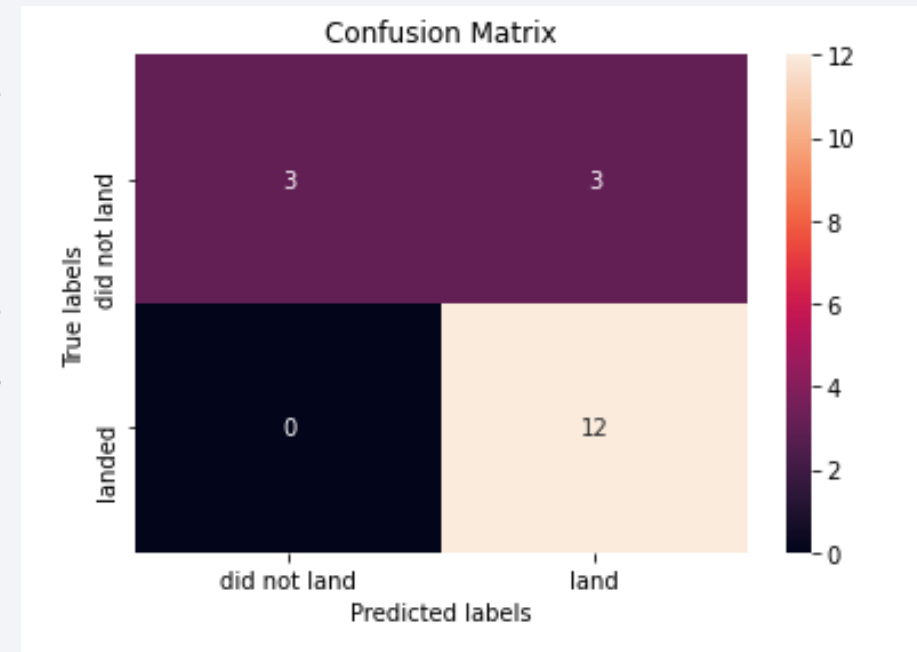
Best model is DecisionTree with a score of 0.8732142857142856

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

# Confusion Matrix

---

- ❖ 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.



# Conclusions

---

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!

