



IBM Developer  
SKILLS NETWORK

# Winning Space Race with Data Science

Niccolò Burberi  
March 2024



# Outline



- Executive Summary [3](#)
  - Introduction [4](#)
    - Methodology [5](#)
    - Insights drawn from EDA [18](#)
      - Launch Sites Proximities Analysis [35](#)
        - Build a Dashboard with Plot Dash [39](#)
          - Predictive Analysis (Classification) [43](#)

# Executive Summary



## Summary of Methodologies

The research attempts to identify the factors for a successful rocket landing. To make this determination, the following methodologies were used:

- Collect data using SpaceX REST API and web scraping techniques
- Wrangle data to create success/fail outcome variable
- Explore data with data visualization techniques, considering the following factors: payload, launch site, flight number and yearly trend
- Analyze the data with SQL, calculating the following statistics: total payload, payload range for successful launches, and total # of successful and failed outcomes
- Explore launch site success rates and proximity to geographical markers
- Visualize the launch sites with the most success and successful payload ranges
- Build Models to predict landing outcomes using logistic regression, support vector machine (SVM), decision tree and K-nearest neighbor (KNN)

## Summary of Results

### Exploratory Data Analysis:

- Launch success has improved over time
- KSC LC-39A has the highest success rate among landing sites
- Orbit ES-L1, GEO, HEO, and SSO have a 100% success rate

### Visualization/Analytics:

- Most launch sites are near the equator, and all are close to the coast

### Predictive Analytics:

- All models performed similarly on the test set. The decision tree model slightly outperformed

# Introduction



SpaceX is a revolutionary company who has disrupted the space industry by offering rocket launches specifically Falcon 9 as low as 62 million dollars; while other providers cost upward of 165 million dollars each. Most of this saving thanks to SpaceX's astounding idea to reuse the first stage of the launch by re-landing the rocket to be used on the next mission. Repeating this process will make the price even lower. As a data scientist of a startup rivaling SpaceX, the goal of this project is to create the machine learning pipeline to predict the landing outcome of the first stage in the future. This project is crucial in identifying the right price to bid against SpaceX for a rocket launch.

The problems included:

- Identifying all factors that influence the landing outcome.
- The relationship between each variables and how it is affecting the outcome.
- The best condition needed to increase the probability of successful landing.



Section 1

# Methodology

# Methodology



- **Collect** data using SpaceX REST API and web scraping techniques
- **Wrangle** data - by filtering the data, handling missing values and applying one hot encoding - to prepare the data for analysis and modeling
- **Explore** data via EDA with SQL and data visualization techniques
- **Visualize** the data using Folium and Plotly Dash
- **Build Models** to predict landing outcomes using classification models. Tune and evaluate models to find best model and parameters

# Data Collection



Data collection is the process of gathering and measuring information on targeted variables in an established system, which then enables one to answer relevant questions and evaluate outcomes. As mentioned, the dataset was collected by REST API and Web Scrapping from Wikipedia.

For REST API, its started by using the get request. Then, we decoded the response content as Json and turn it into a pandas dataframe using `json_normalize()`. We then cleaned the data, checked for missing values and fill with whatever needed.

For web scrapping, we will use the BeautifulSoup to extract the launch records as HTML table, parse the table and convert it to a pandas dataframe for further analysis.

# Data Collection SpaceX API



- 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()`

From:

[https://github.com/niccolo1994/Applied-Data-Science-Capstone/blob/2b1ec07c6f09d3dd3f81b4790ee6391f05ae286f/01\\_jupyter-labs-spacex-data-collection-api.ipynb](https://github.com/niccolo1994/Applied-Data-Science-Capstone/blob/2b1ec07c6f09d3dd3f81b4790ee6391f05ae286f/01_jupyter-labs-spacex-data-collection-api.ipynb)

```
spacex_url="https://api.spacexdata.com/v4/launches/past"

response = requests.get(spacex_url)

# Use json_normalize method to convert the json result into a dataframe
data = pd.json_normalize(response.json())

# Lets take a subset of our dataframe keeping only the features we want and the flight number, and date_utc.
data = data[['rocket', 'payloads', 'launchpad', 'cores', 'flight_number', 'date_utc']]

# We will remove rows with multiple cores because those are falcon rockets with 2 extra rocket boosters and rows
data = data[data['cores'].map(len)==1]
data = data[data['payloads'].map(len)==1]

# Since payloads and cores are lists of size 1 we will also extract the single value in the list and replace the
data['cores'] = data['cores'].map(lambda x : x[0])
data['payloads'] = data['payloads'].map(lambda x : x[0])

# We also want to convert the date_utc to a datetime datatype and then extracting the date leaving the time
data['date'] = pd.to_datetime(data['date_utc']).dt.date

# Using the date we will restrict the dates of the launches
data = data[data['date'] <= datetime.date(2020, 11, 13)]

# Create a data from launch_dict
data = pd.DataFrame(launch_dict)

# Hint data[ 'BoosterVersion' ]!= 'Falcon 1'
data_falcon9 = data[data.BoosterVersion == 'Falcon 9']
data_falcon9

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

# Data Collection Scraping



- 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

From:

[https://github.com/niccolo1994/Applied-Data-Science-Capstone/blob/2b1ec07c6f09d3dd3f81b4790ee6391f05ae286f/03\\_jupyter-labs-webscraping.ipynb](https://github.com/niccolo1994/Applied-Data-Science-Capstone/blob/2b1ec07c6f09d3dd3f81b4790ee6391f05ae286f/03_jupyter-labs-webscraping.ipynb)

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

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

```
# 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")
print(html_tables)
```

```
launch_dict= dict.fromkeys(column_names)

# Remove an irrelevant column
del launch_dict['Date and time ( )']

# Let's initial the launch_dict with each value to be an empty list
launch_dict['Flight No.'] = []
launch_dict['Launch site'] = []
launch_dict['Payload'] = []
launch_dict['Payload mass'] = []
launch_dict['Orbit'] = []
launch_dict['Customer'] = []
launch_dict['Launch outcome'] = []
# Added some new columns
launch_dict['Version Booster']=[]
launch_dict['Booster landing']=[]
launch_dict['Date']=[]
launch_dict['Time']=[]
```

```
df = pd.DataFrame.from_dict(launch_dict)
df.head()
```

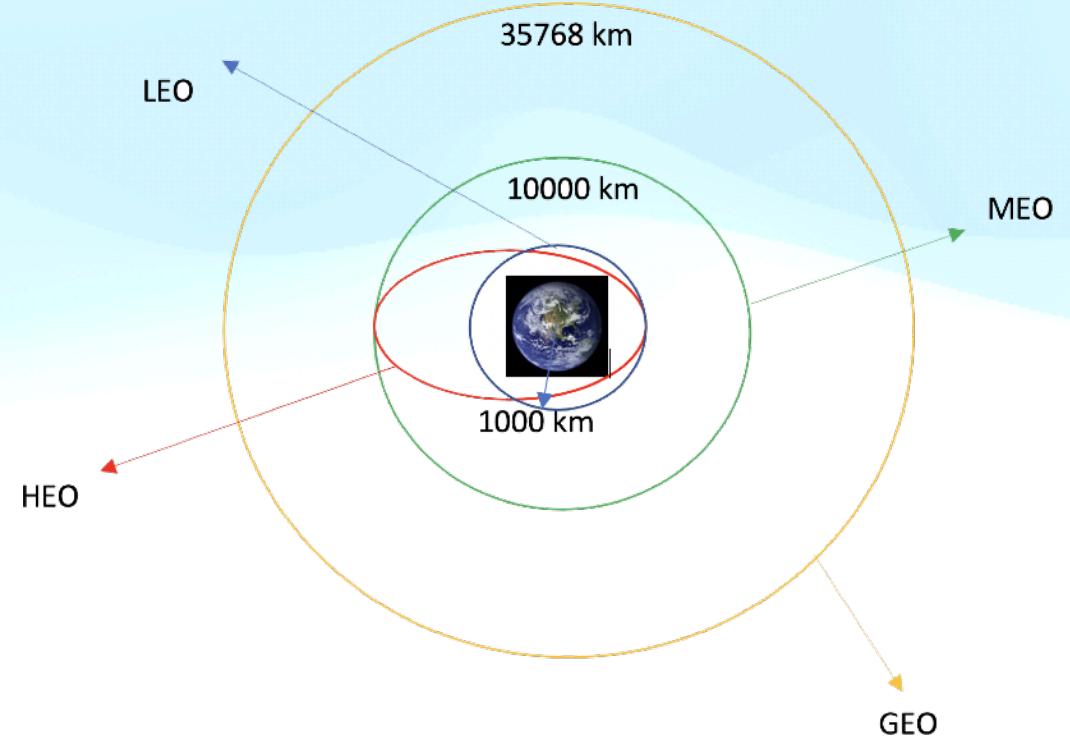
# Data Wrangling



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

From:

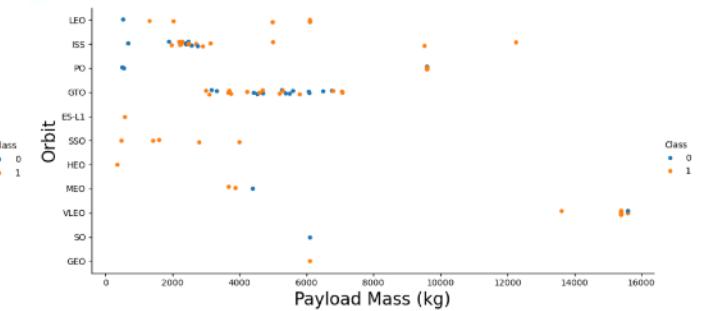
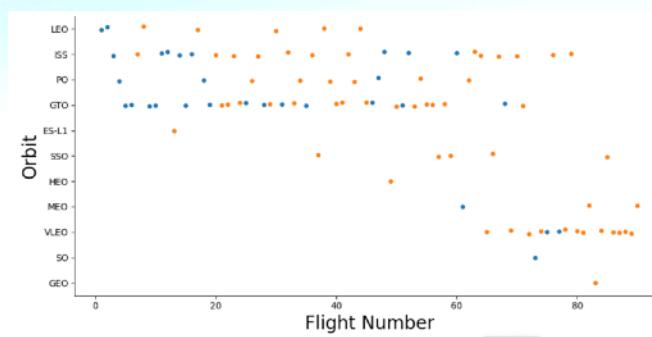
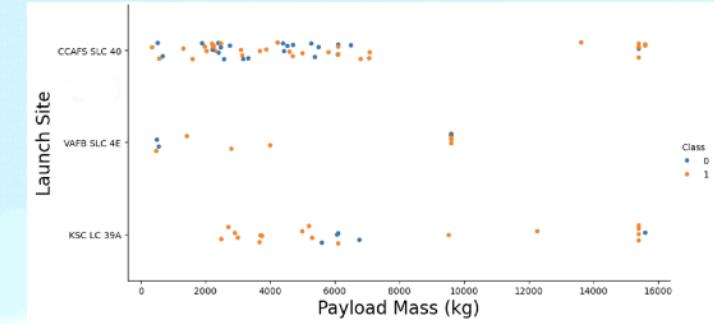
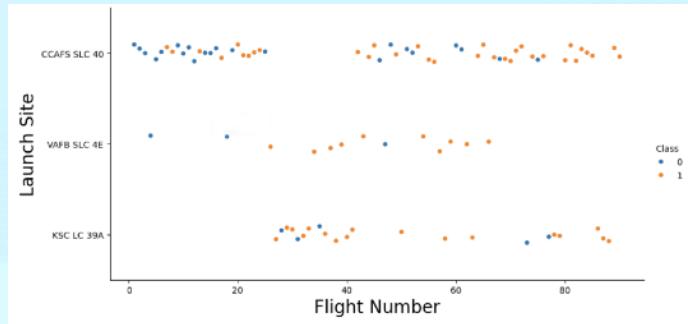
[https://github.com/niccolo1994/Applied-Data-Science-Capstone/blob/2b1ec07c6f09d3dd3f81b4790ee6391f05ae286f/02\\_labs-jupyter-spacex-Data%20wrangling.ipynb](https://github.com/niccolo1994/Applied-Data-Science-Capstone/blob/2b1ec07c6f09d3dd3f81b4790ee6391f05ae286f/02_labs-jupyter-spacex-Data%20wrangling.ipynb)



# EDA with Data Visualization



- View relationship by using scatter plots. The variables could be useful for machine learning if a relationship exists:
  - Flight Number and Launch Site.
  - Payload and Launch Site.
  - Flight Number and Orbit Type.
  - Payload and Orbit Type.
- Show comparisons among discrete categories with bar charts. Bar charts show the relationships among the categories and a measured value.



From:

[https://github.com/niccolo1994/Applied-Data-Science-Capstone/blob/2b1ec07c6f09d3dd3f81b4790ee6391f05ae286f/05\\_jupyter-labs-eda-dataviz.ipynb.jupyterlite.ipynb](https://github.com/niccolo1994/Applied-Data-Science-Capstone/blob/2b1ec07c6f09d3dd3f81b4790ee6391f05ae286f/05_jupyter-labs-eda-dataviz.ipynb.jupyterlite.ipynb)

# EDA with Data Visualization



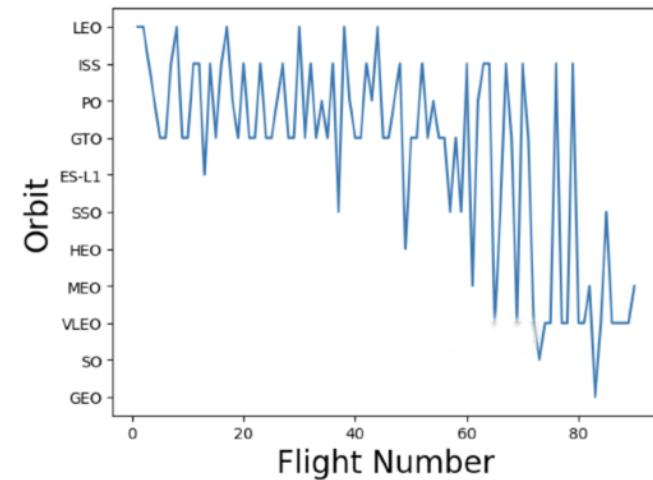
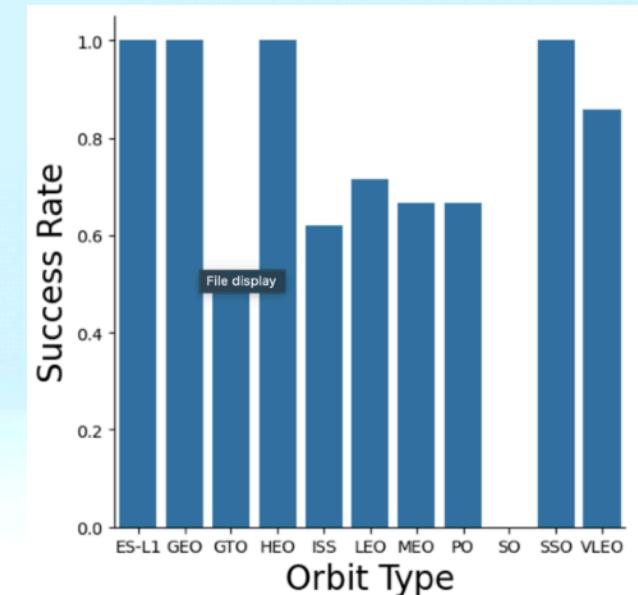
We will use scatter plot visualization tools such as bar graph and line plots graph for further analysis.

Bar graphs is one of the easiest way to interpret the relationship between the attributes. In this case, we will use the bar graph to determine which orbits have the highest probability of success.

We then use the line graph to show a trends or pattern of the attribute over time which in this case, is used for see the relationship between orbits and flight number.

From:

[https://github.com/niccolo1994/Applied-Data-Science-Capstone/blob/2b1ec07c6f09d3dd3f81b4790ee6391f05ae286f/05\\_jupyter-labs-eda-dataviz.ipynb.jupyterlite.ipynb](https://github.com/niccolo1994/Applied-Data-Science-Capstone/blob/2b1ec07c6f09d3dd3f81b4790ee6391f05ae286f/05_jupyter-labs-eda-dataviz.ipynb.jupyterlite.ipynb)



# EDA with SQL



Using SQL, we had performed many queries to get better understanding of the dataset, Ex:

Display:

- Names of unique launch sites
- 5 records where launch site begins with 'CCA'
- Total payload mass carried by boosters launched by NASA (CRS)
- Average payload mass carried by booster version F9 v1.1.

List:

- Date of first successful landing on ground pad
- Names of boosters which had success landing on drone ship and have payload mass greater than 4,000 but less than 6,000
- Total number of successful and failed missions
- Names of booster versions which have carried the max payload
- Failed landing outcomes on drone ship, their booster version and launch site for the months in the year 2015
- Count of landing outcomes between 2010-06-04 and 2017-03-20 (desc)

From:

[https://github.com/niccolo1994/Applied-Data-Science-Capstone/blob/2b1ec07c6f09d3dd3f81b4790ee6391f05ae286f/04\\_jupyter-labs-eda-sql-coursera\\_sqlite.ipynb](https://github.com/niccolo1994/Applied-Data-Science-Capstone/blob/2b1ec07c6f09d3dd3f81b4790ee6391f05ae286f/04_jupyter-labs-eda-sql-coursera_sqlite.ipynb)

# Build an Interactive Map with Folium



## Markers Indicating Launch Sites

- Added blue circle at NASA Johnson Space Center's coordinate with a popup label showing its name using its latitude and longitude coordinates
- Added red circles at all launch sites coordinates with a popup label showing its name using its name using its latitude and longitude coordinates Map with Folium

## Colored Markers of Launch Outcomes

- Added colored markers of successful (green) and unsuccessful (red) launches at each launch site to show which launch sites have high success rates

## Distances Between a Launch Site to Proximities

- Added colored lines to show distance between launch site CCAFS SLC-40 and its proximity to the nearest coastline, railway, highway, and city

From:

[https://github.com/niccolo1994/Applied-Data-Science-Capstone/blob/2b1ec07c6f09d3dd3f81b4790ee6391f05ae286f/06\\_lab\\_jupyter\\_launch\\_site\\_location.jupyterlite.ipynb](https://github.com/niccolo1994/Applied-Data-Science-Capstone/blob/2b1ec07c6f09d3dd3f81b4790ee6391f05ae286f/06_lab_jupyter_launch_site_location.jupyterlite.ipynb)

# Build a Dashboard with Plotly Dash



We built an interactive dashboard with Plotly dash which allowing the user to play around with the data as they need.

We plotted pie charts showing the total launches by a certain sites.

We then plotted scatter graph showing the relationship with Outcome and Payload Mass (Kg) for the different booster version.

From:

[https://github.com/niccolo1994/Applied-Data-Science-Capstone/blob/2b1ec07c6f09d3dd3f81b4790ee6391f05ae286f/  
07\\_Build%20an%20Interactive%20Dashboard%20with%20Ploty%20Dash](https://github.com/niccolo1994/Applied-Data-Science-Capstone/blob/2b1ec07c6f09d3dd3f81b4790ee6391f05ae286f/07_Build%20an%20Interactive%20Dashboard%20with%20Ploty%20Dash)

# Predictive Analysis (Classification)



## Charts

- Create NumPy array from the Class column
- Standardize the data with StandardScaler. Fit and transform the data
- Split the data using train\_test\_split
- Create a GridSearchCV object with cv=10 for parameter optimization
- Apply GridSearchCV on different algorithms: logistic regression (LogisticRegression()), support vector machine (SVC()), decision tree (DecisionTreeClassifier()), K-Nearest Neighbor (KNeighborsClassifier())
- Calculate accuracy on the test data using .score() for all models
- Assess the confusion matrix for all models
- Identify the best model using Jaccard\_Score, F1\_Score and Accuracy

From:

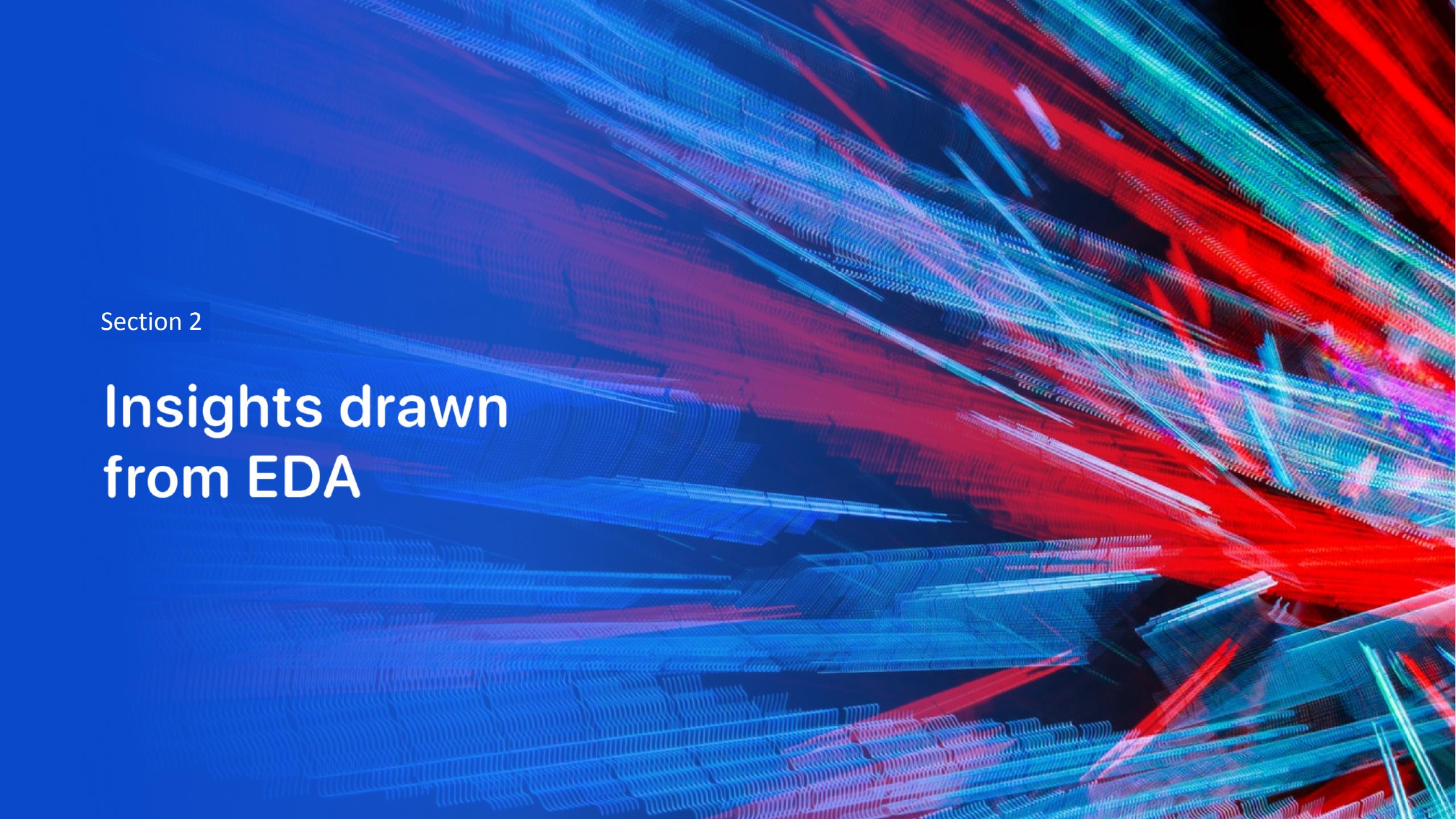
[https://github.com/niccolo1994/Applied-Data-Science-Capstone/blob/2b1ec07c6f09d3dd3f81b4790ee6391f05ae286f/  
08\\_SpaceX\\_Machine\\_Learning\\_Prediction\\_Part\\_5.ipynb](https://github.com/niccolo1994/Applied-Data-Science-Capstone/blob/2b1ec07c6f09d3dd3f81b4790ee6391f05ae286f/08_SpaceX_Machine_Learning_Prediction_Part_5.ipynb)

# Results



The results will be categorized to 3 main results which is:

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

The background of the slide features a complex, abstract digital pattern. It consists of numerous thin, glowing lines that create a sense of depth and motion. The colors used are primarily shades of blue, red, and purple, which are bright against a dark, almost black, background. These lines form a grid-like structure that is more dense and vibrant towards the right side of the frame, while appearing more sparse and blurred towards the left.

Section 2

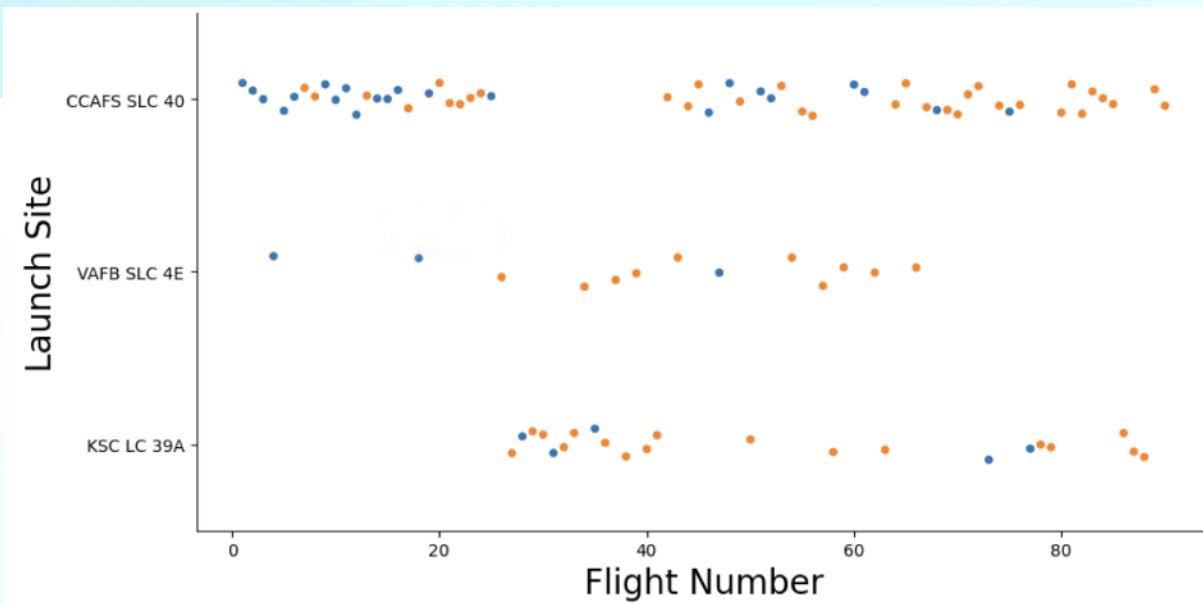
## Insights drawn from EDA

# Flight Number vs. Launch Site



This scatter plot shows that the larger the flights amount of the launch site, the greater the success rate will be.

However, site CCAFS SLC40 shows the least pattern of this.

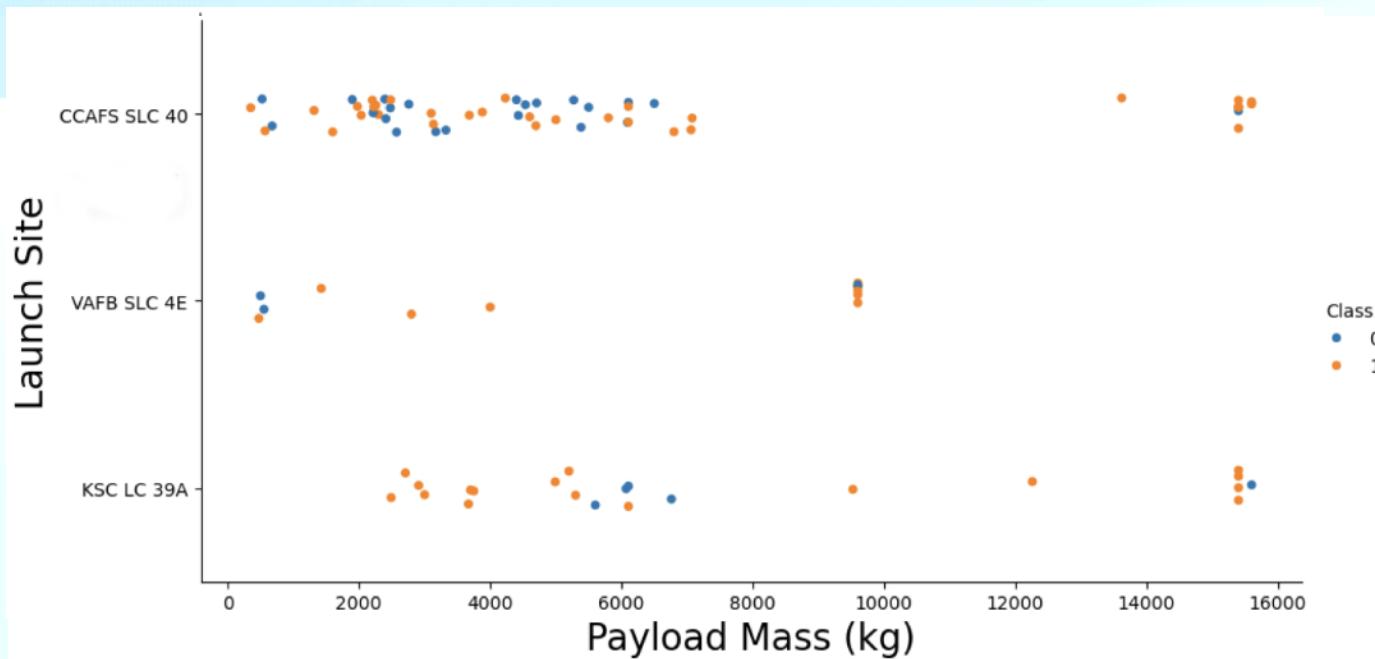


# Payload vs. Launch Site



This scatter plot shows once the payload mass is greater than 7000kg, the probability of the success rate will be highly increased.

However, there is no clear pattern to say the launch site is dependent to the payload mass for the success rate.

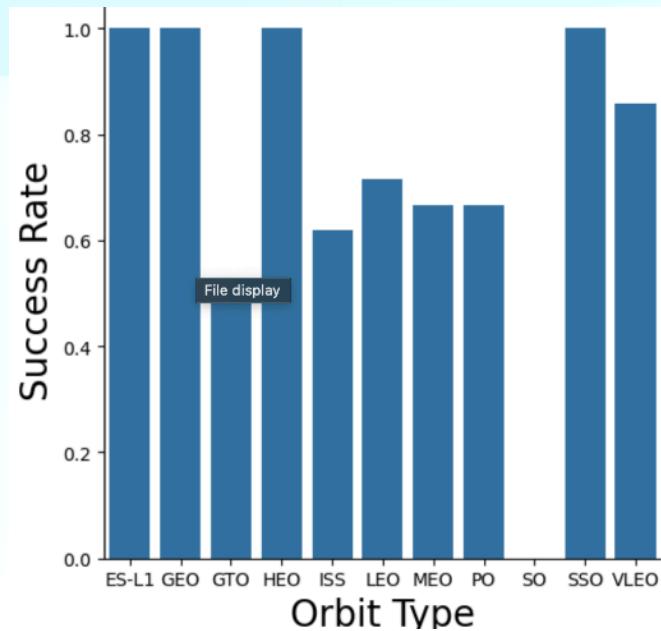


# Success Rate vs. Orbit Type

SPACEX

This scatter plot shows once the pay load mass is greater than 7000kg, the probability of the success rate will be highly increased.

However, there is no clear pattern to say the launch site is dependent to the pay load mass for the success rate.

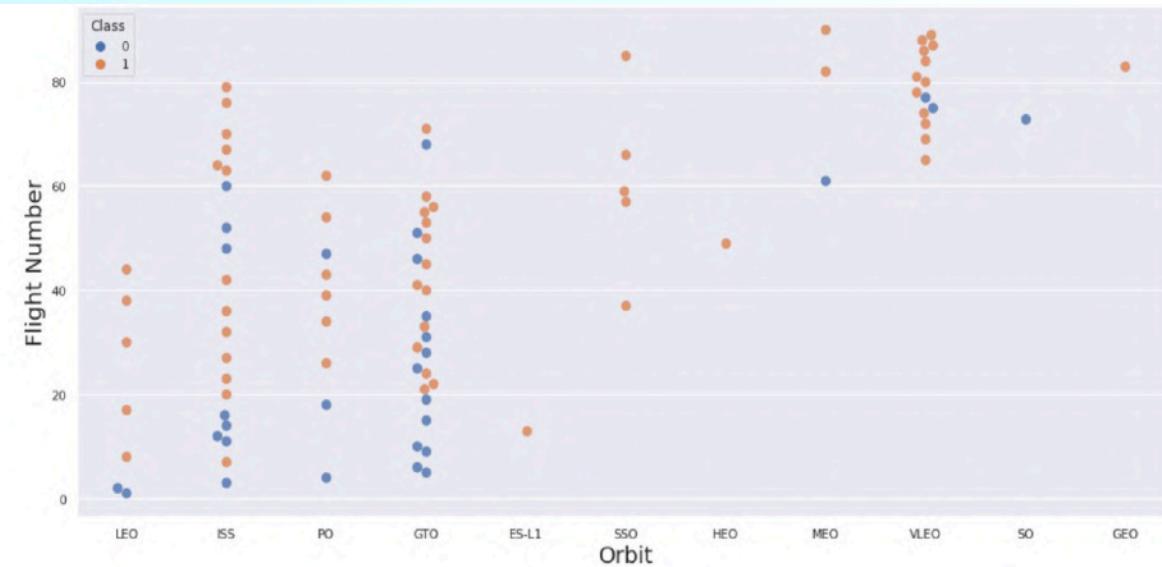


# Flight Number vs. Orbit Type

SPACEX

This scatter plot shows that generally, the larger the flight number on each orbits, the greater the success rate (especially LEO orbit) except for GTO orbit which depicts no relationship between both attributes.

Orbit that only has 1 occurrence should also be excluded from above statement as it's needed more dataset.



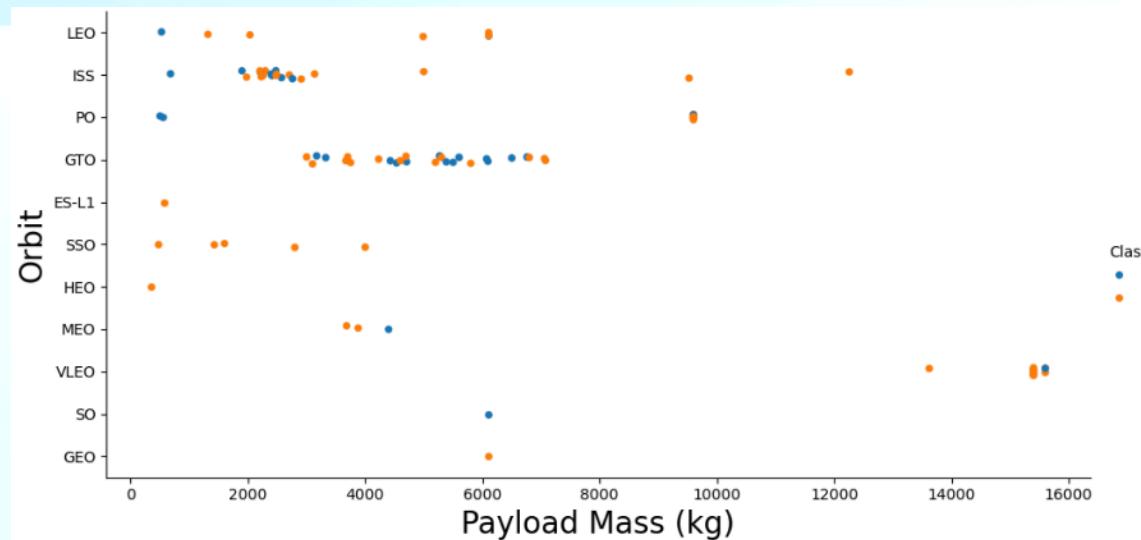
# Payload vs. Orbit Type



Heavier payload has positive impact on LEO, ISS and PO orbit. However, it has negative impact on MEO and VLEO orbit.

GTO orbit seem to depict no relation between the attributes.

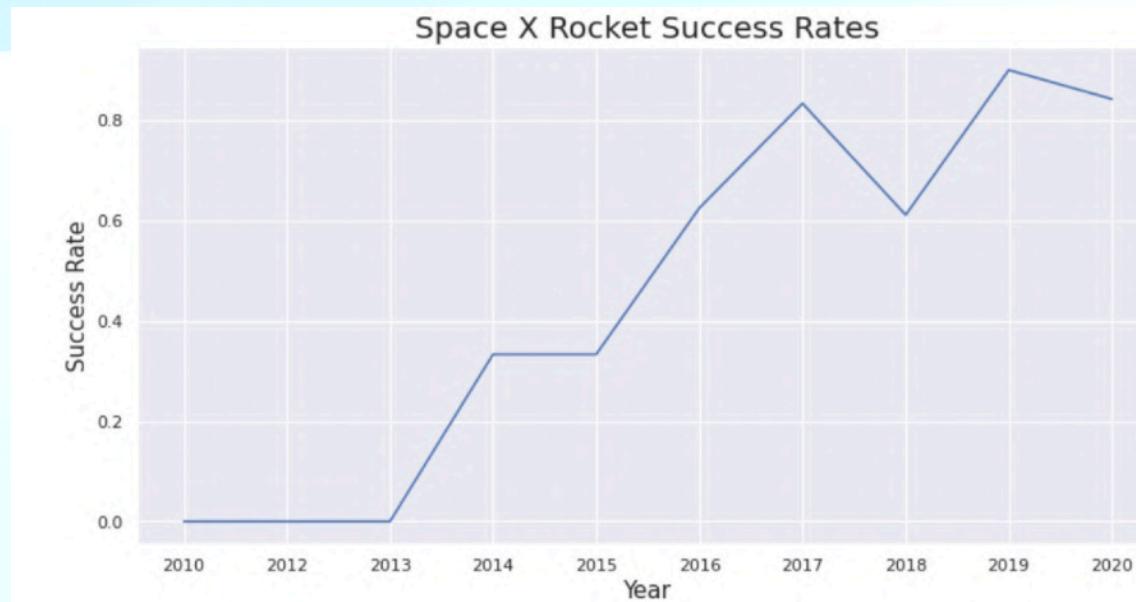
Meanwhile, again, SO, GEO and HEO orbit need more dataset to see any pattern or trend.



# Launch Success Yearly Trend



This figure clearly depicted an increasing trend from the year 2013 until 2020.  
If this trend continues for the next year onward, the success rate will steadily increase until reaching 100% success rate.



# All Launch Site Names



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

In [5]:

```
%sql SELECT DISTINCT LAUNCH_SITE as "Launch_Sites" FROM SPACEX;  
  
* ibm_db_sa://zpw86771:***@fbdb88901-ebdb-4a4f-a32e-9822b9fb237b.clogj3  
sd0tgtu0lqde00.databases.appdomain.cloud:32731/bludb  
Done.
```

Out[5]:

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

# Launch Site Names Begin with 'CCA'

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

```
In [11]: task_2 = """
    SELECT *
    FROM SpaceX
    WHERE LaunchSite LIKE 'CCA%'
    LIMIT 5
"""
create_pandas_df(task_2, database=conn)
```

	date	time	boosterversion	launchsite	payload	payloadmasskg	orbit	customer	missionoutcome	landingoutcome
0	2010-04-06	18:45:00	F9 v1.0 B0003	CCAFS LC-40	Dragon Spacecraft Qualification Unit	0	LEO	SpaceX	Success	Failure (parachute)
1	2010-08-12	15:43:00	F9 v1.0 B0004	CCAFS LC-40	Dragon demo flight C1, two CubeSats, barrel of...	0	LEO (ISS)	NASA (COTS) NRO	Success	Failure (parachute)
2	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
3	2012-08-10	00:35:00	F9 v1.0 B0006	CCAFS LC-40	SpaceX CRS-1	500	LEO (ISS)	NASA (CRS)	Success	No attempt
4	2013-01-01	00:00:00	F9 v1.0 B0007	CCAFS LC-40	Dragon CRS-2	---	LEO	SpaceX	Success	Success

# Total Payload Mass



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

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

```
%sql SELECT SUM(PAYLOAD_MASS__KG_) AS "Total Payload Mass by NASA (CRS)
```

```
* ibm_db_sa://zpw86771:***@fdb88901-ebdb-4a4f-a32e-9822b9fb237b.c1ogj3
sd0tgtu0lqde00.databases.appdomain.cloud:32731/bludb
Done.
```

**Total Payload Mass by NASA (CRS)**

---

45596

# Average Payload Mass by F9 v1.1



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

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

```
%sql SELECT AVG(PAYLOAD_MASS__KG_) AS "Average Payload Mass by Booster  
WHERE BOOSTER_VERSION = 'F9 v1.1';
```

```
* ibm_db_sa://zpw86771:***@fbdb88901-ebdb-4a4f-a32e-9822b9fb237b.clogj3  
sd0tgtu0lqde00.databases.appdomain.cloud:32731/bludb  
Done.
```

**Average Payload Mass by Booster Version F9 v1.1**

---

2928

# First Successful Ground Landing Date



We use the min() function to find the result

We observed that the dates of the first successful landing outcome on ground pad was  
22nd December 2015

```
%sql SELECT MIN(DATE) AS "First Succesful Landing Outcome in Ground Pad"  
WHERE LANDING_OUTCOME = 'Success (ground pad)';
```

```
* ibm_db_sa://zpw86771:***@fbdb88901-ebdb-4a4f-a32e-9822b9fb237b.clogj3  
sd0tgtu0lqde00.databases.appdomain.cloud:32731/bludb  
Done.
```

**First Succesful Landing Outcome in Ground Pad**

---

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

```
%sql SELECT BOOSTER_VERSION FROM SPACEX WHERE LANDING_OUTCOME = 'Success (drone ship)' \
AND PAYLOAD_MASS_KG_ > 4000 AND PAYLOAD_MASS_KG_ < 6000;

* ibm_db_sa://zpw86771:****@fdb88901-ebdb-4a4f-a32e-9822b9fb237b.clogj3sd0tgtu0lgde00.databases.appdomain.cloud:32731/bludb
Done.

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 MissionOutcome was a success or a failure.

List the total number of successful and failure mission outcomes

```
%sql SELECT COUNT(MISSION_OUTCOME) AS "Successful Mission" FROM SPACEX WHERE MISSION_OUTCOME LIKE 'Success%';
```

```
* ibm_db_sa://zpw86771:***@fb88901-ebdb-4a4f-a32e-9822b9fb237b.clogj3sd0tgtu0lqde00.databases.appdomain.cloud:32731/bludb  
Done.
```

**Successful Mission**

100

```
%sql SELECT COUNT(MISSION_OUTCOME) AS "Failure Mission" FROM SPACEX WHERE MISSION_OUTCOME LIKE 'Failure%';
```

```
* ibm_db_sa://zpw86771:***@fb88901-ebdb-4a4f-a32e-9822b9fb237b.clogj3sd0tgtu0lqde00.databases.appdomain.cloud:32731/bludb  
Done.
```

**Failure Mission**

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.

```
sql SELECT DISTINCT BOOSTER_VERSION AS "Booster Versions which carried the Maximum Payload Mass" FROM SPACEX  
WHERE PAYLOAD_MASS_KG_ =(SELECT MAX(PAYLOAD_MASS_KG_) FROM SPACEX);
```

```
* ibm_db_sa://zpw86771:****@fdb88901-ebdb-4a4f-a32e-9822b9fb237b.clogj3sd0tgtu0lgde00.databases.appdomain.cloud:32731/bludb
```

```
Done.
```

**Booster Versions which carried the Maximum Payload Mass**

F9 B5 B1048.4

F9 B5 B1048.5

F9 B5 B1049.4

F9 B5 B1049.5

F9 B5 B1049.7

F9 B5 B1051.3

F9 B5 B1051.4

F9 B5 B1051.6

F9 B5 B1056.4

F9 B5 B1058.3

F9 B5 B1060.2

F9 B5 B1060.3

# 2015 Launch Records



We used a combination 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

```
%sql SELECT BOOSTER_VERSION, LAUNCH_SITE FROM SPACEX WHERE DATE LIKE '2015-%' AND \
LANDING_OUTCOME = 'Failure (drone ship)';

* ibm_db_sa://zpw86771:***@fdb88901-ebdb-4a4f-a32e-9822b9fb237b.c1ogj3sd0tgtu0lqde00.
databases.appdomain.cloud:32731/bludb
Done.

booster_version    launch_site
F9 v1.1 B1012    CCAFS LC-40
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 applied the GROUP BY clause to group the landing outcomes and the ORDER BY clause to order the grouped landing outcome in descending order.

```
%sql SELECT LANDING_OUTCOME as "Landing Outcome", COUNT(LANDING_OUTCOME) AS "Total Count" FROM SPACEX \
WHERE DATE BETWEEN '2010-06-04' AND '2017-03-20' \
GROUP BY LANDING_OUTCOME \
ORDER BY COUNT(LANDING_OUTCOME) DESC ;
```

```
* ibm_db_sa://zpw86771:****@fdb88901-ebdb-4a4f-a32e-9822b9fb237b.clogj3sd0tgtu01qde00.databases.appdomain.c
loud:32731/bludb
Done.
```

Landing Outcome	Total 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

The background of the slide is a photograph taken from space at night. It shows the curvature of the Earth against a dark blue-black void of space. City lights are visible as numerous small white and yellow dots, primarily concentrated in the lower right quadrant where the United States appears. In the upper left quadrant, the green and blue glow of the aurora borealis is visible in the upper atmosphere.

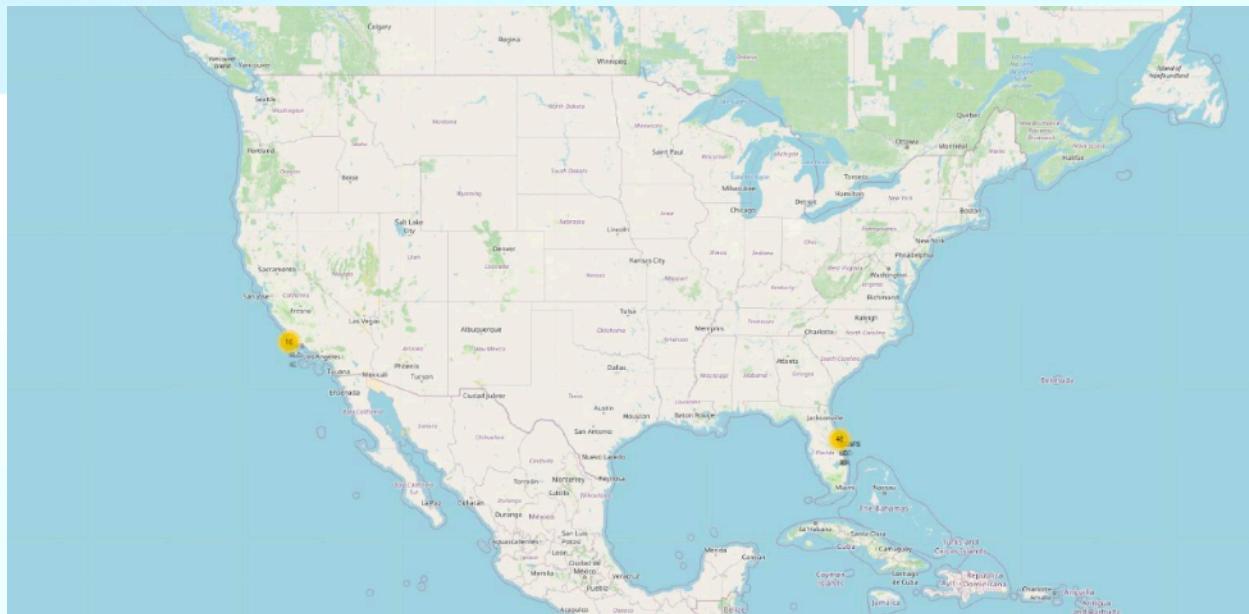
Section 3

# Launch Sites Proximities Analysis

# Launch Sites

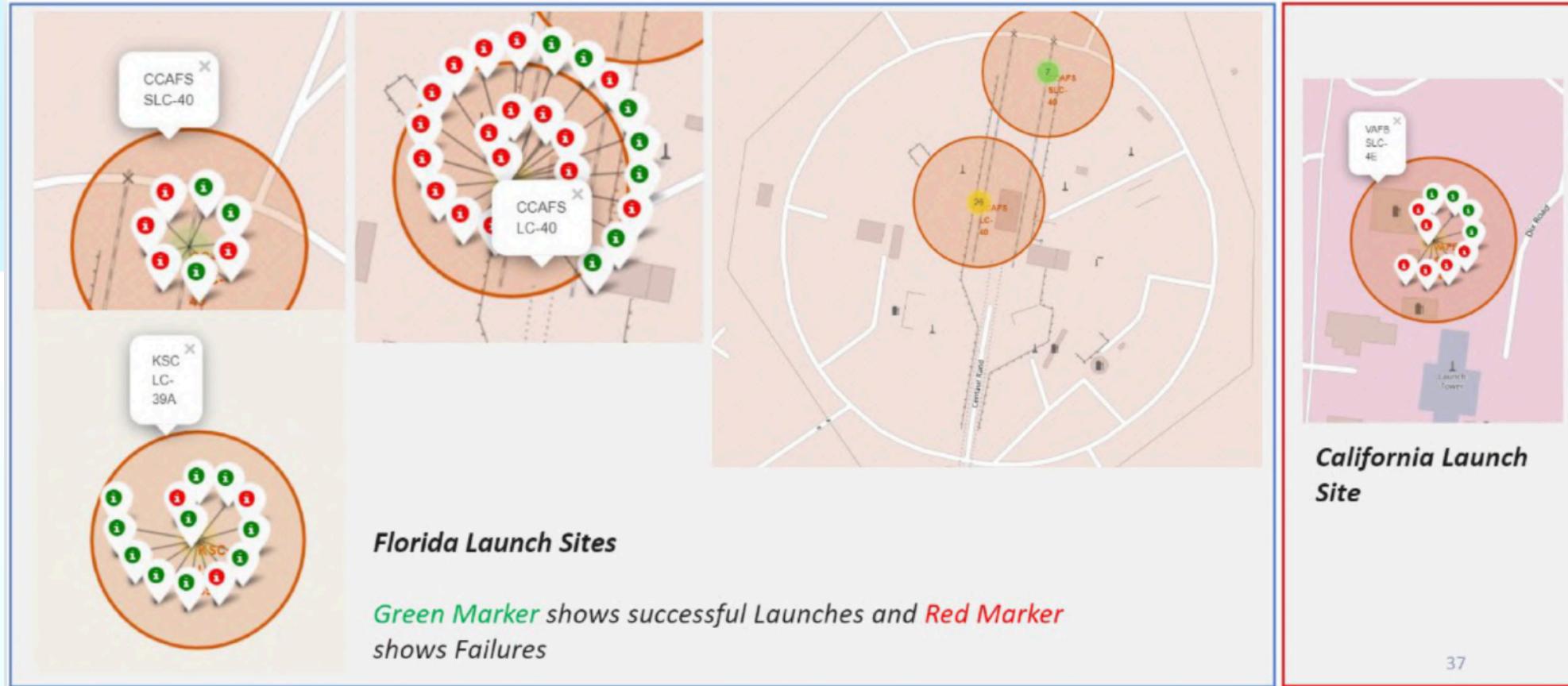
SPACEX

Near Equator: the closer the launch site to the equator, the easier it is to launch to equatorial orbit, and the more help you get from Earth's rotation for a prograde orbit. Rockets launched from sites near the equator get an additional natural boost- due to the rotational speed of earth - that helps save the cost of putting in extra fuel and boosters.

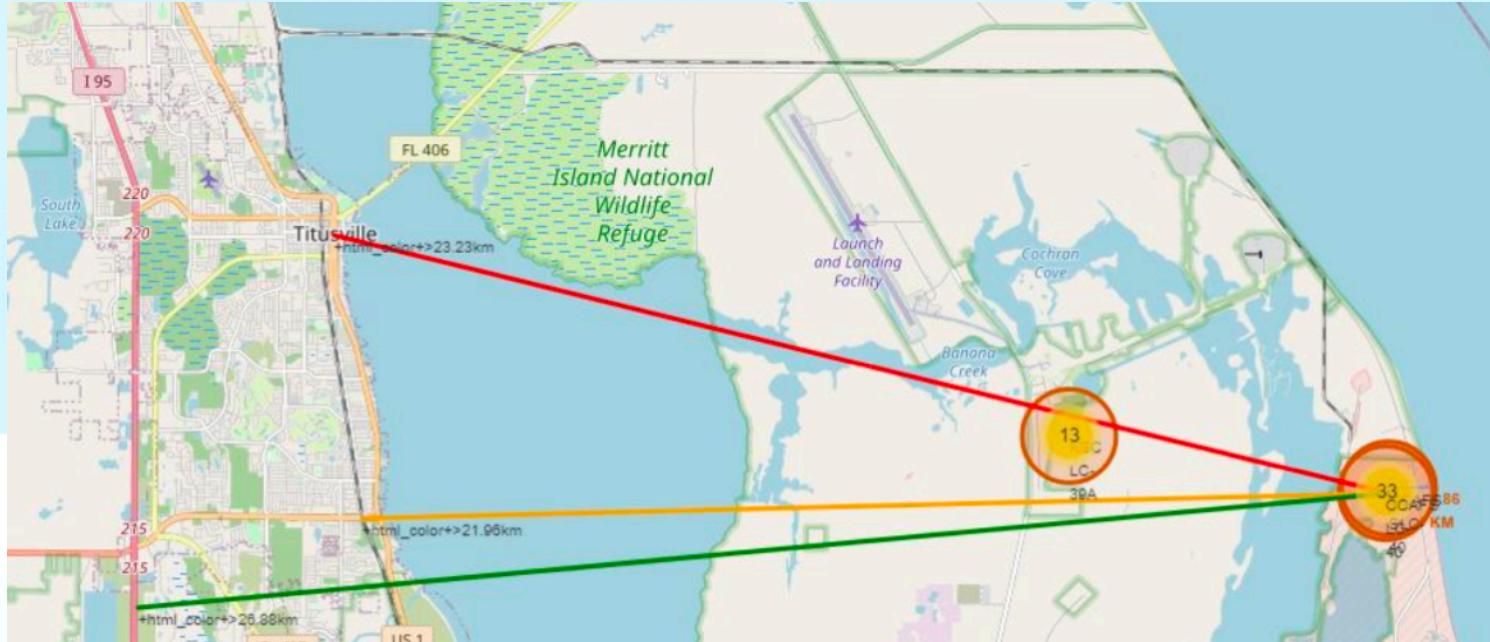


# Launch Outcomes

SPACEX



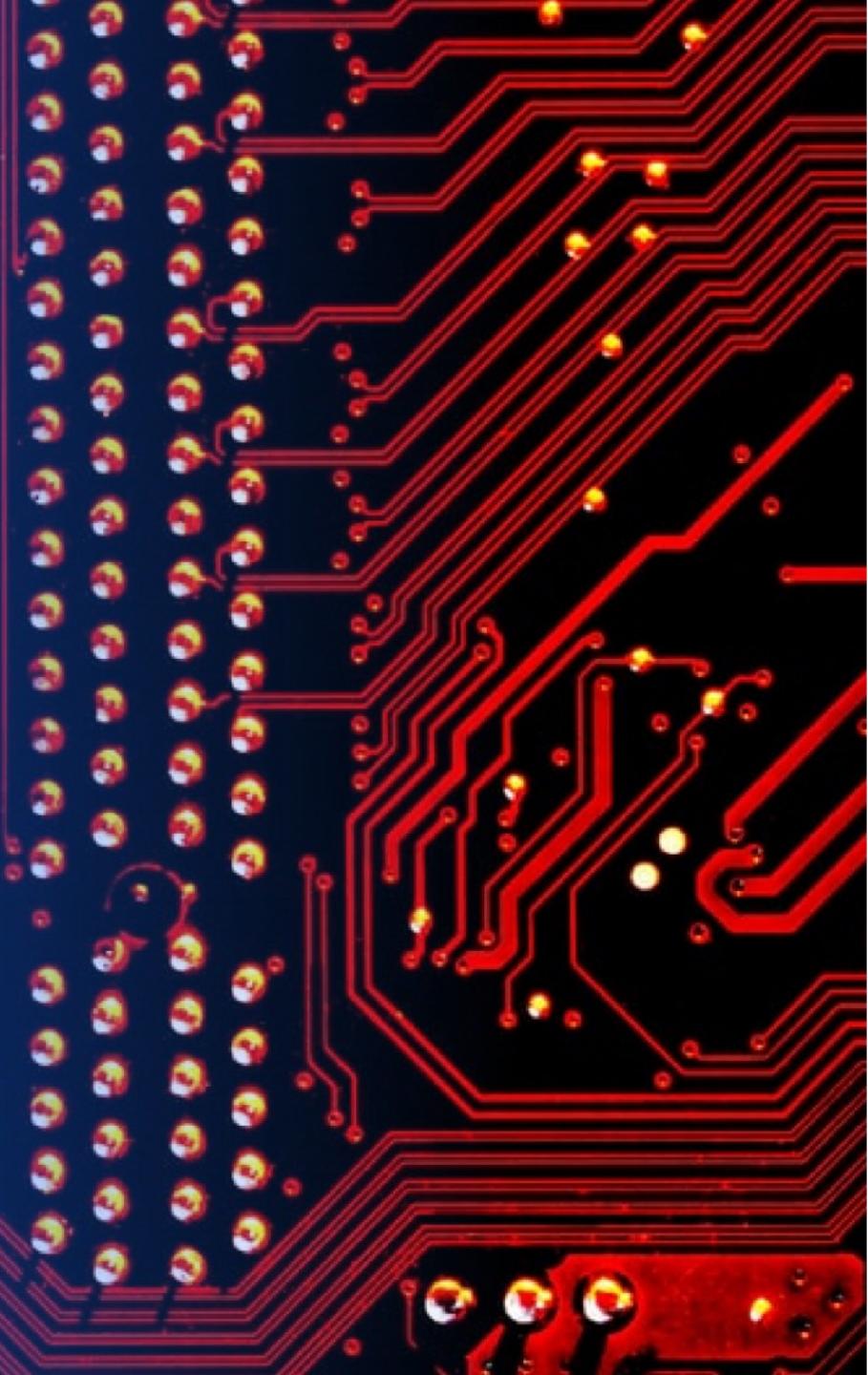
# Launch Sites Distance to Landmarks



- Are launch sites in close proximity to railways? **NO**
- Are launch sites in close proximity to highways? **NO**
- Are launch sites in close proximity to coastline? **YES**
- Do launch sites keep certain distance away from cities? **YES**

Section 4

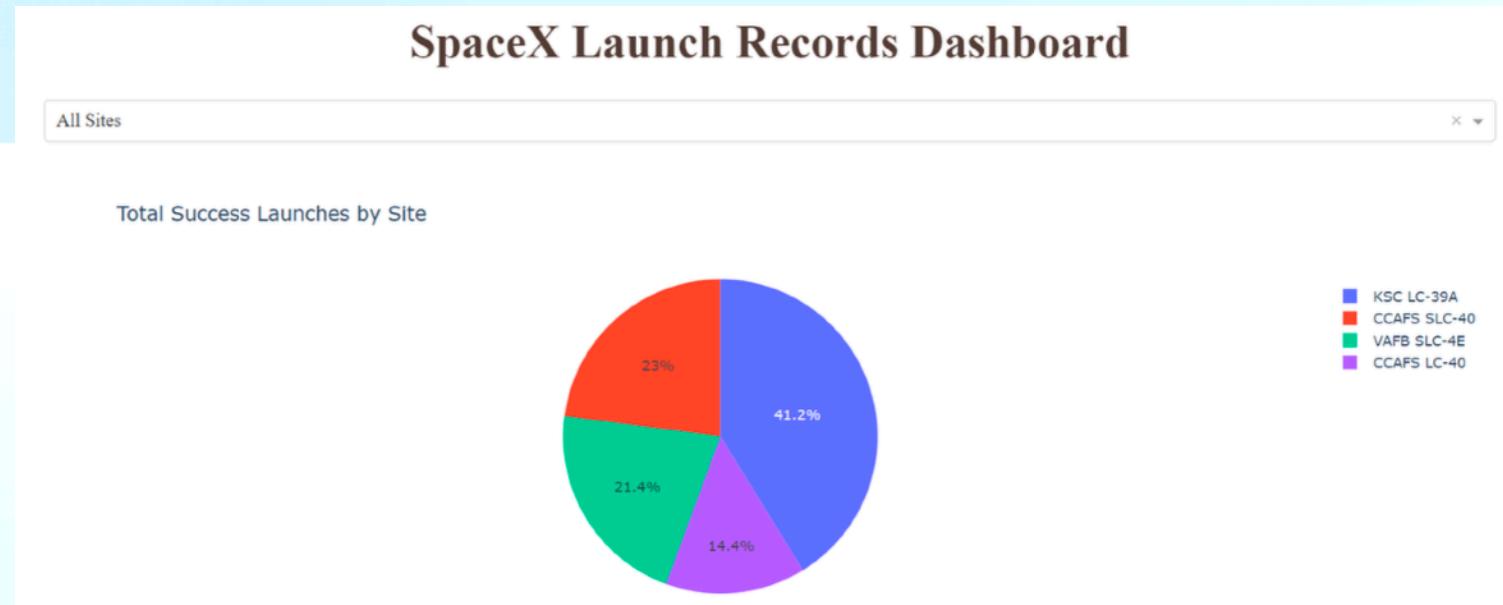
# Build a Dashboard with Plotly Dash



# Launch Success by Site



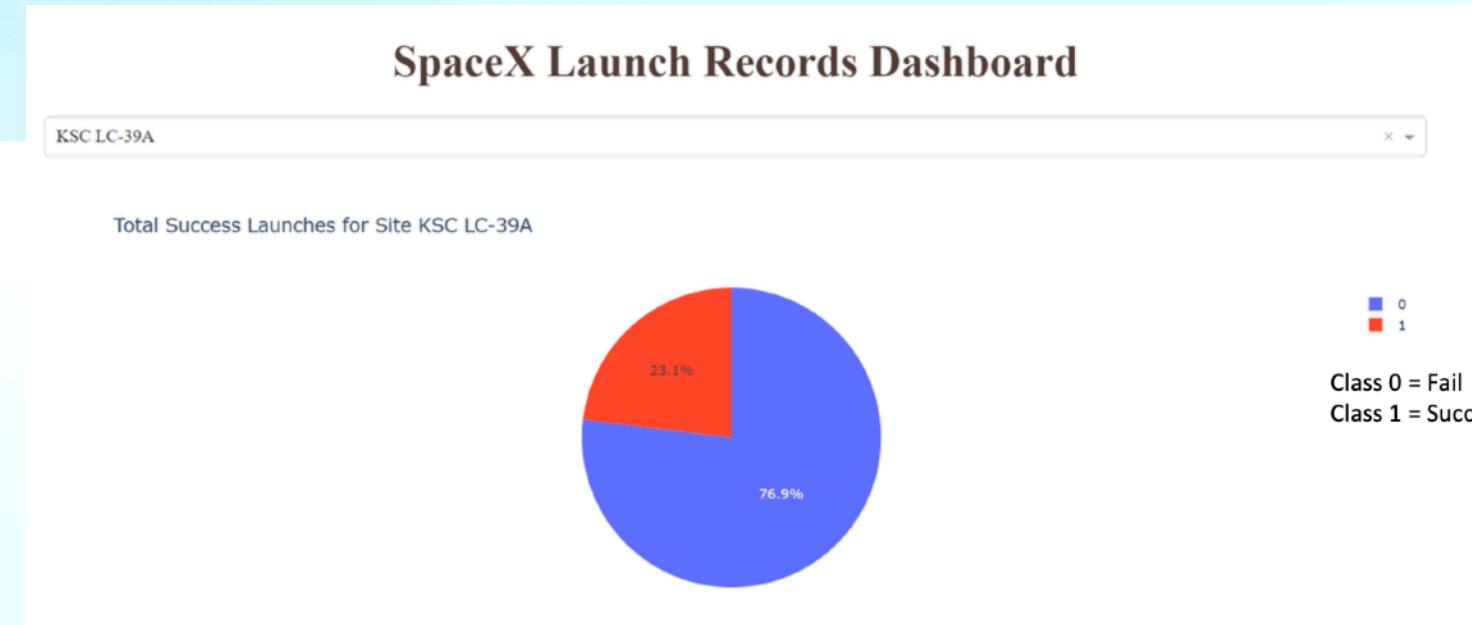
Success as Percent of Total: **KSC LC-39A** has the most successful launches amongst launch sites (**41.2%**)



# Launch Success (KSC LC-29A)



Success as Percent of Total: **KSC LC-39A** has the highest success rate amongst launch sites (**76.9%**)



# Payload Mass vs. Success vs. Booster Version Category

SPACEX

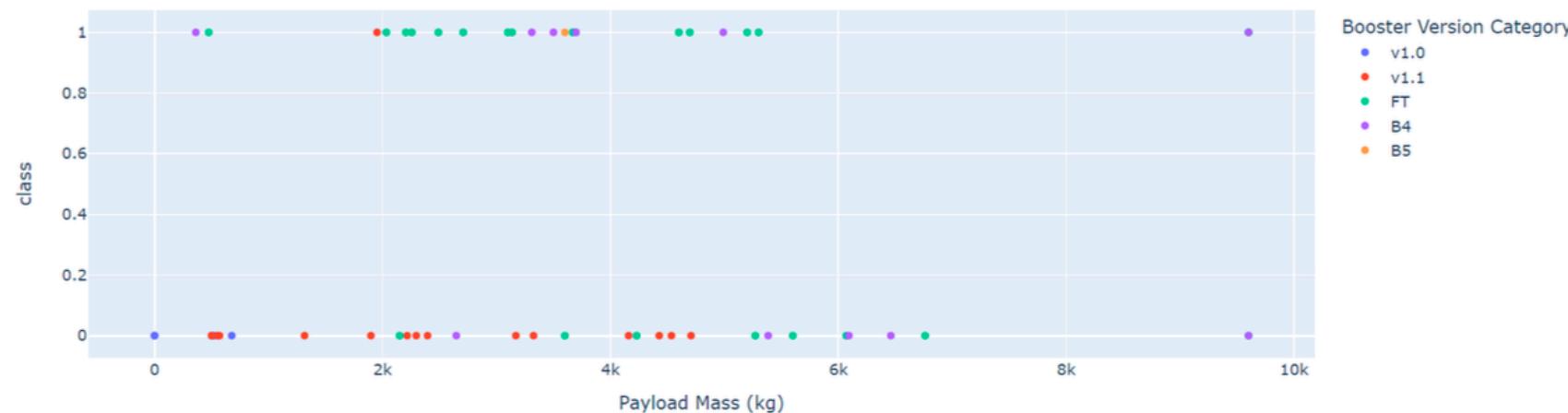
## By Booster Version

- Payloads between 2,000 kg and 5,000 kg have the highest success rate
- 1 indicating successful outcome and 0 indicating an unsuccessful outcome

Payload range (Kg):



Correlation Between Payload and Success for All Sites



The background of the slide features a dynamic, abstract design. It consists of several thick, curved lines that transition from a bright yellow at the top right to a deep blue at the bottom left. These lines create a sense of motion and depth, resembling a tunnel or a stylized landscape. The overall effect is modern and professional.

Section 5

# Predictive Analysis (Classification)

# Classification Accuracy



- All the models performed at about the same level and had the same scores and accuracy. This is likely due to the small dataset. The Decision Tree model slightly outperformed the rest when looking at `.best_score_`

```
algorithms = {'KNN':knn_cv.best_score_, 'Tree':tree_cv.best_score_, 'LogisticRegression':logreg_cv.best_score_}
bestalgorithm = max(algorithms, key=algorithms.get)
print('Best Algorithm is',bestalgorithm,'with a score of',algorithms[bestalgorithm])
if bestalgorithm == 'Tree':
    print('Best Params is :',tree_cv.best_params_)
if bestalgorithm == 'KNN':
    print('Best Params is :',knn_cv.best_params_)
if bestalgorithm == 'LogisticRegression':
    print('Best Params is :',logreg_cv.best_params_)

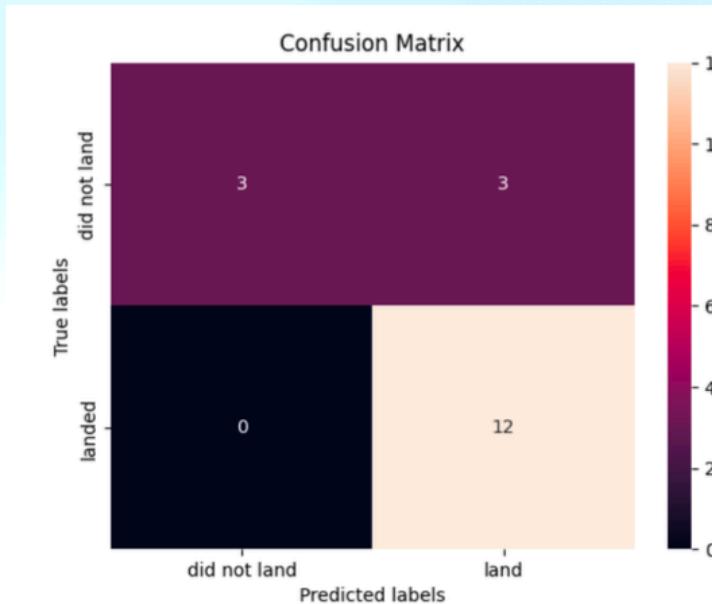
Best Algorithm is Tree with a score of 0.9017857142857142
Best Params is : {'criterion': 'entropy', 'max_depth': 10, 'max_features': 'auto', 'min_samples_leaf': 2, 'min_samples_split': 10, 'splitter': 'random'}
```

# 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.

- Precision =  $TP / (TP + FP)$
- $12 / 15 = .80$
- Recall =  $TP / (TP + FN)$
- $12 / 12 = 1$
- F1 Score =  
 $2 * (Precision * Recall) / (Precision + Recall)$
- $2 * (.8 * 1) / (.8 + 1) = .89$
- Accuracy =  
 $(TP + TN) / (TP + TN + FP + FN) = .833$



# Conclusions



- **Model Performance:** The models performed similarly on the test set with the decision tree model slightly outperforming
- **Equator:** Most of the launch sites are near the equator for an additional natural boost - due to the rotational speed of earth - which helps save the cost of putting in extra fuel and boosters
- **Coast:** All the launch sites are close to the coast
- **Launch Success:** Increases over time
- **KSC LC-39A:** Has the highest success rate among launch sites. Has a 100% success rate for launches less than 5,500 kg
- **Orbits:** ES-L1, GEO, HEO, and SSO have a 100% success rate
- **Payload Mass:** Across all launch sites, the higher the payload mass (kg), the higher the success rate

# Appendix



GitHub: <https://github.com/niccolo1994/Applied-Data-Science-Capstone.git>

Thank you!

