

# Winning Space Race with Data Science

Muhammad Umair  
Salim

28/04/2023

[This is the Link](#) for Github repo of this  
project



# Outline

---

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

# Executive Summary

---

## **Methodology**

Collected data from SpaceX API and by web scraping Wikipedia pages. After collecting data cleaned in and performed exploratory data analysis using SQL pandas and visualization libraries. Using the insights drawn by EDA further explored the data by interactive visualization and creating interactive dashboards for finding results. Finally completed the task of building a predictive classification model.

## **Results**

Results were drawn from these three steps:

- EDA
- Interactive Visualization
- Predictive Analysis

After drawing results they were then compiled into this report for purpose of presentation.

# Introduction

---

- SpaceX launches Falcon 9 rockets at a cost of around \$62m. This is considerably cheaper than other providers (which usually cost upwards of \$165m), and much of the savings are because SpaceX can land, and then re-use the first stage of the rocket.
- If we can make predictions on whether the first stage will land, we can determine the cost of a launch, and use this information to assess whether or not an alternate company should bid and SpaceX for a rocket launch.
- This project will ultimately predict if the Space X Falcon 9 first stage will land successfully.





Section 1

# Methodology

# Methodology Summery

---

## 1. Data Collection

- Making GET requests to the SpaceX REST API
- Web Scraping

## 2. Data Wrangling

- Using the `.fillna()` method to remove NaN values
- Using the `.value_counts()` method to determine the following:
  - Number of launches on each site
  - Number and occurrence of each orbit
  - Number and occurrence of mission outcome per orbit type
- Creating a landing outcome label that shows the following:
  - 0 when the booster did not land successfully
  - 1 when the booster did land successfully

## 3. Exploratory Data Analysis

- Using SQL queries to manipulate and evaluate the SpaceX dataset
- Using Pandas and Matplotlib to visualize relationships between variables, and determine patterns

## 4. Interactive Visual Analytics

- Geospatial analytics using Folium
- Creating an interactive dashboard using Plotly Dash

## 5. Data Modelling and Evaluation

- Using Scikit-Learn to:
  - Pre-process (standardize) the data
  - Split the data into training and testing data using `train_test_split`
  - Train different classification models
  - Find hyperparameters using GridSearchCV
- Plotting confusion matrices for each classification model
- Assessing the accuracy of each classification model

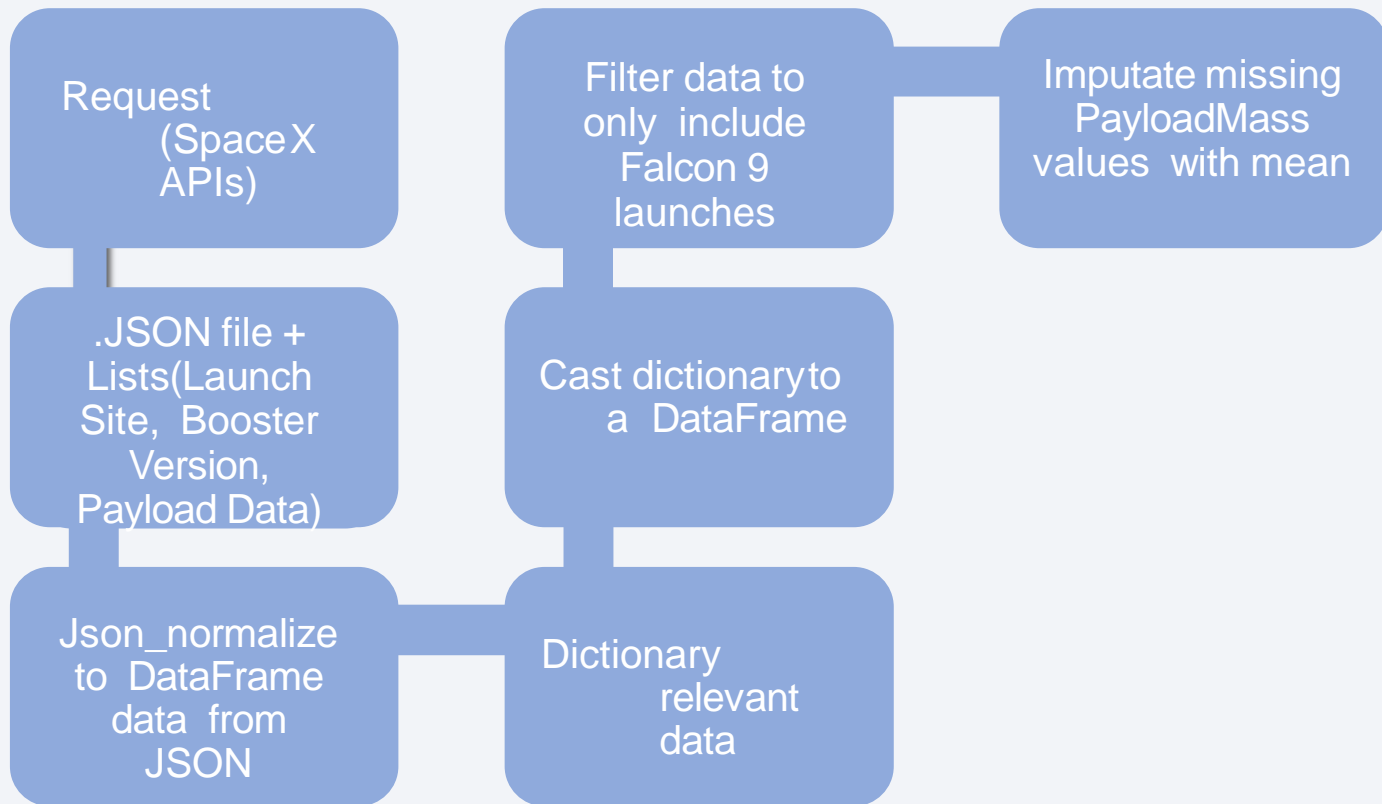
## 4. Interactive Visual Analytics

- Geospatial analytics using Folium
- Creating an interactive dashboard using Plotly Dash

## 5. Data Modelling and Evaluation

- Using Scikit-Learn to:
  - Pre-process (standardize) the data
  - Split the data into training and testing data using `train_test_split`
  - Train different classification models
  - Find hyperparameters using GridSearchCV
- Plotting confusion matrices for each classification model
- Assessing the accuracy of each classification model

# Data Collection – SpaceX API



[Github Link](#)

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

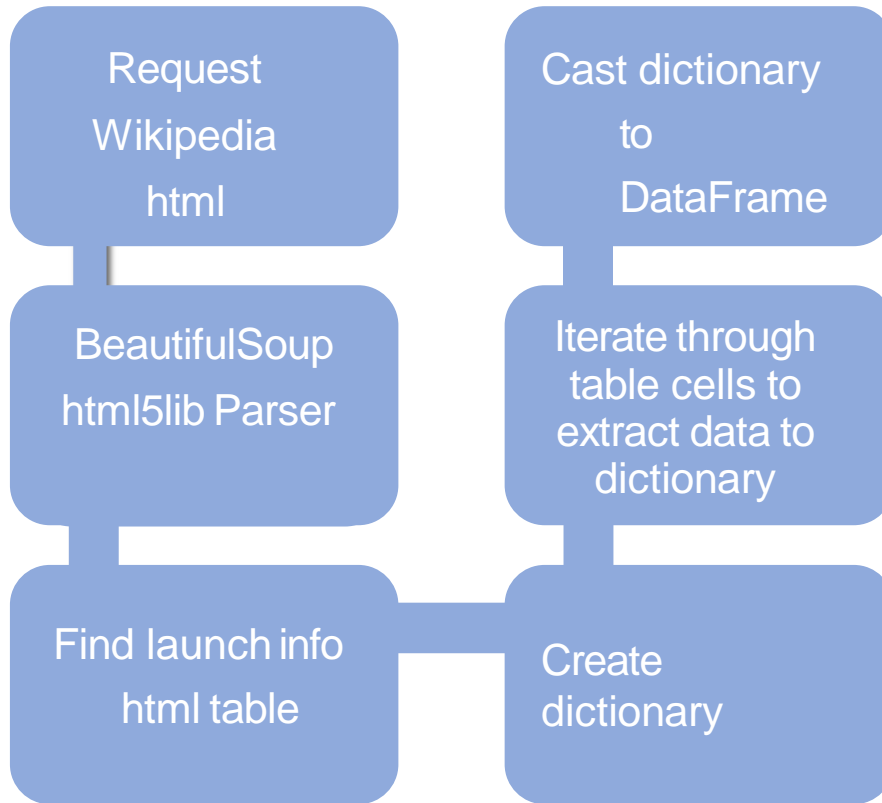
In [7]: response = requests.get(spacex_url)
```

```
In [10]: response.status_code

Out[10]: 200
```

```
In [12]: # Use json_normalize meethod to convert the json result into a dataframe
data = pd.json_normalize(response.json())
```

# Data Collection – Scraping



[Github Link](#)

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
        response = requests.get(static_url)
```

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

```
In [7]: # Use BeautifulSoup() to create a BeautifulSoup object from a response text content
        soup = BeautifulSoup(response.content, "html.parser")
```

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

```
In [8]: # Use soup.title attribute
        print(soup.title)

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



# Data Wrangling

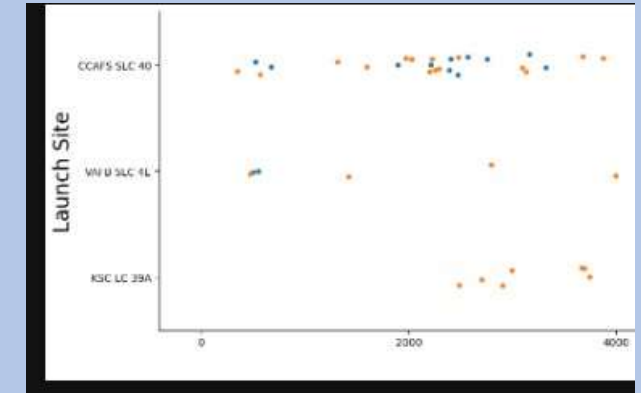
---

- Create a training label with landing outcomes where successful = 1 & failure = 0.
  - Outcome column has two components: 'Mission Outcome' 'Landing Location'
  - New training label column 'class' with a value of 1 if 'Mission Outcome' is True and 0 otherwise. Value Mapping:
  - True ASDS, True RTLS, & True Ocean – set to -> 1
  - None None, False ASDS, None ASDS, False Ocean, False RTLS – set to -> 0
- [Github Link](#)

# EDA with Data Visualization

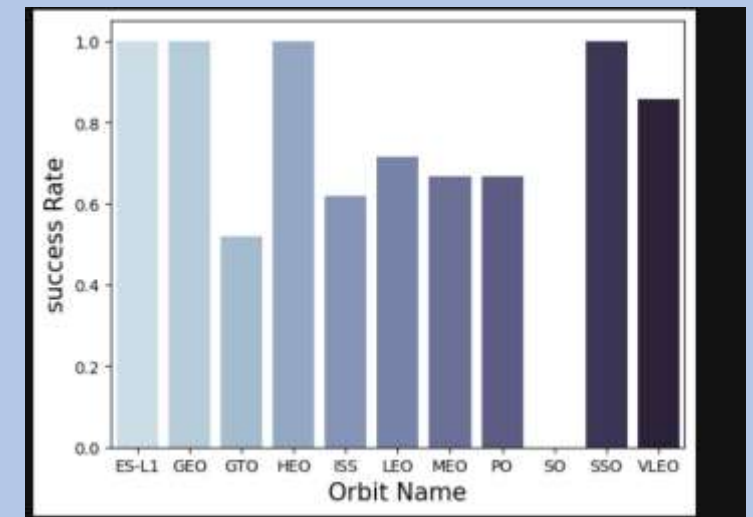
- Exploratory Data Analysis performed on variables Flight Number, Payload Mass, Launch Site, Orbit, Class and Year.
- Plots Used:
- Flight Number vs. Payload Mass, Flight Number vs. Launch Site, Payload Mass vs. Launch Site, Orbit vs. Success Rate, Flight Number vs. Orbit, Payload vs Orbit, and Success Yearly Trend
- Scatter plots, line charts, and bar plots were used to compare relationships between variables to
- decide if a relationship exists so that they could be used in training the machine learning model

[Github Link](#)



```
# Hint: Use get_dummies() function on the categorical columns
features_one_hot = pd.get_dummies(features, columns=['Orbit', 'LaunchSite', 'LandingPad', 'Serial'])
features_one_hot
```

FlightNumber	PayloadMass	Flights	GridFins	Reused	Legs	Block	\
1	6104.958412	1	False	False	False	1.0	
2	525.000000	1	False	False	False	1.0	
3	677.000000	1	False	False	False	1.0	



# EDA with SQL

- Loaded data set into IBM DB2 Database.
- Queried using SQL Python integration.
- Queries were made to get a better understanding of the dataset.
- Queried information about launch site names, mission outcomes, various pay load sizes of customers and booster versions, and landing outcomes

[Github Link](#)

```
Display the names of the unique launch sites in the s

In [6]: %sql
SELECT DISTINCT LAUNCH_SITE FROM SPACEX

* ibm_db_sa://nwc66894:***@824dfd4d-99de-440d-99
Done.

Out[6]: launch_site
        CCAFS LC-40
        CCAFS SLC-40
        KSC LC-39A
        VAFB SLC-4E
```

```
List the total number of successful and failure mission outcomes

In [23]: %sql
SELECT MISSION_OUTCOME, COUNT(*) AS num_missions
FROM SPACEX
GROUP BY MISSION_OUTCOME;

* ibm_db_sa://nwc66894:***@824dfd4d-99de-440d-9991-629c01b38
Done.

Out[23]: mission_outcome  num_missions
        Failure (in flight)      1
        Success                  99
        Success (payload status unclear)  1
```

# Build an Interactive Map with Folium

---

- The following steps were taken to visualize the launch data on an interactive map:
  1. Mark all launch sites on a map
    - Initialise the map using a Folium `Map` object
    - Add a `folium.Circle` and `folium.Marker` for each launch site on the launch map
  2. Mark the success/failed launches for each site on a map
    - As many launches have the same coordinates, it makes sense to cluster them together.
    - Before clustering them, assign a marker colour of successful (class = 1) as green, and failed (class = 0) as red.
    - To put the launches into clusters, for each launch, add a `folium.Marker` to the `MarkerCluster()` object.
    - Create an icon as a text label, assigning the `icon_color` as the `marker_colour` determined previously.
  3. Calculate the distances between a launch site to its proximities
    - To explore the proximities of launch sites, calculations of distances between points can be made using the `Lat` and `Long` values.
    - After marking a point using the `Lat` and `Long` values, create a `folium.Marker` object to show the distance.
    - To display the distance line between two points, draw a `folium.PolyLine` and add this to the map.

[Github Link](#)

# Build a Dashboard with Plotly Dash

---

- The following plots were added to a Plotly Dash dashboard to have an interactive visualisation of the data:
1. Pie chart (`px.pie()`) showing the total successful launches per site
    - This makes it clear to see which sites are most successful
    - The chart could also be filtered (using a `dcc.Dropdown()` object) to see the success/failure ratio for an individual site
  2. Scatter graph (`px.scatter()`) to show the correlation between outcome (success or not) and payload mass (kg)
    - This could be filtered (using a `RangeSlider()` object) by ranges of payload masses
    - It could also be filtered by booster version

[Github Link](#)



# Predictive Analysis (Classification)

---

## Model Development



- To prepare the dataset for model development:
  - Load dataset
  - Perform necessary data transformations (standardise and pre-process)
  - Split data into training and test data sets, using `train_test_split()`
  - Decide which type of machine learning algorithms are most appropriate
- For each chosen algorithm:
  - Create a `GridSearchCV` object and a dictionary of parameters
  - Fit the object to the parameters
  - Use the training data set to train the model

## Model Evaluation



- For each chosen algorithm:
  - Using the output `GridSearchCV` object:
    - Check the tuned hyperparameters (`best_params_`)
    - Check the accuracy (`score` and `best_score_`)
  - Plot and examine the Confusion Matrix

## Finding the Best Classification Model



- Review the accuracy scores for all chosen algorithms
- The model with the highest accuracy score is determined as the best performing model

[Github Link](#)



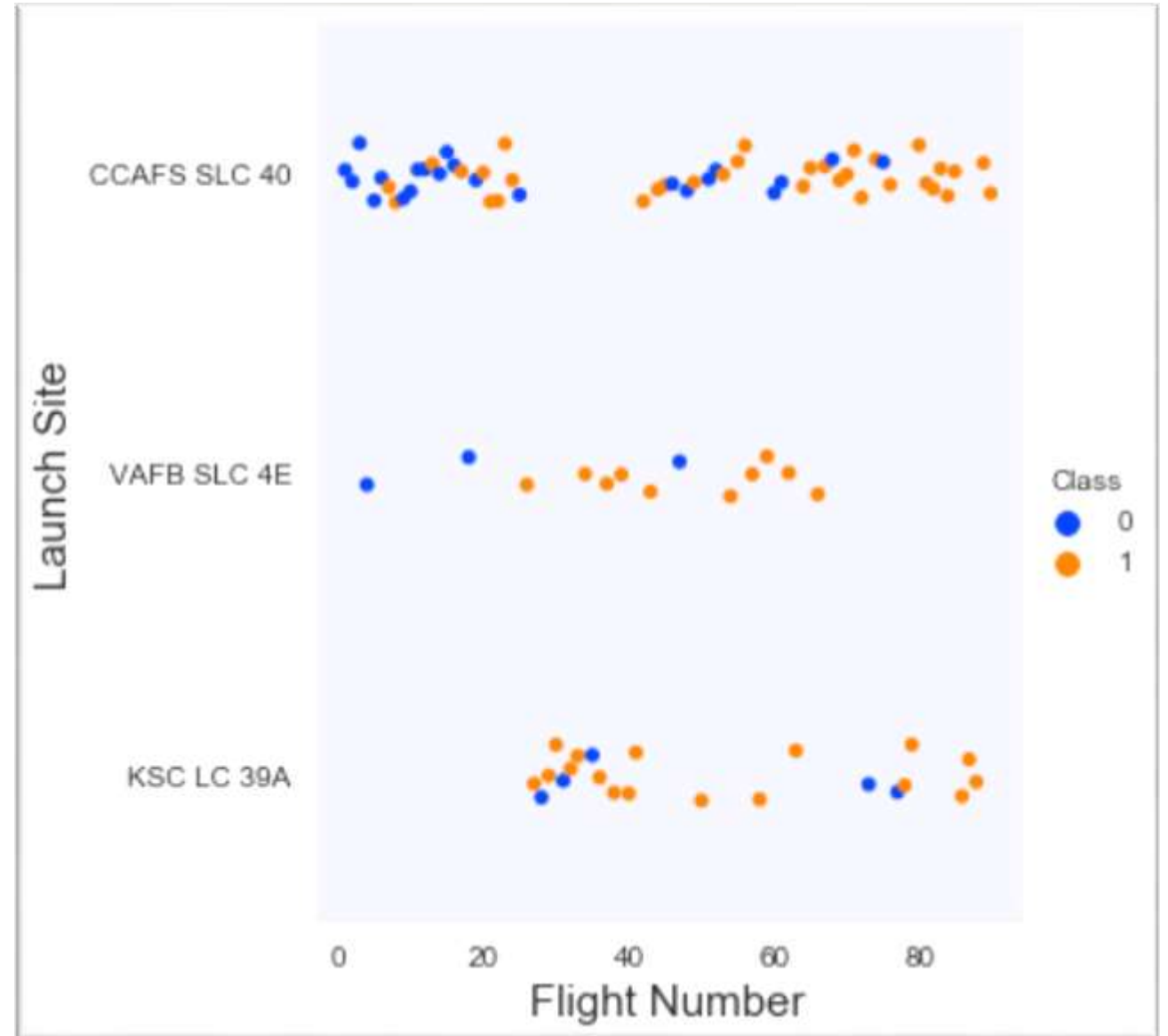
# **Results**

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

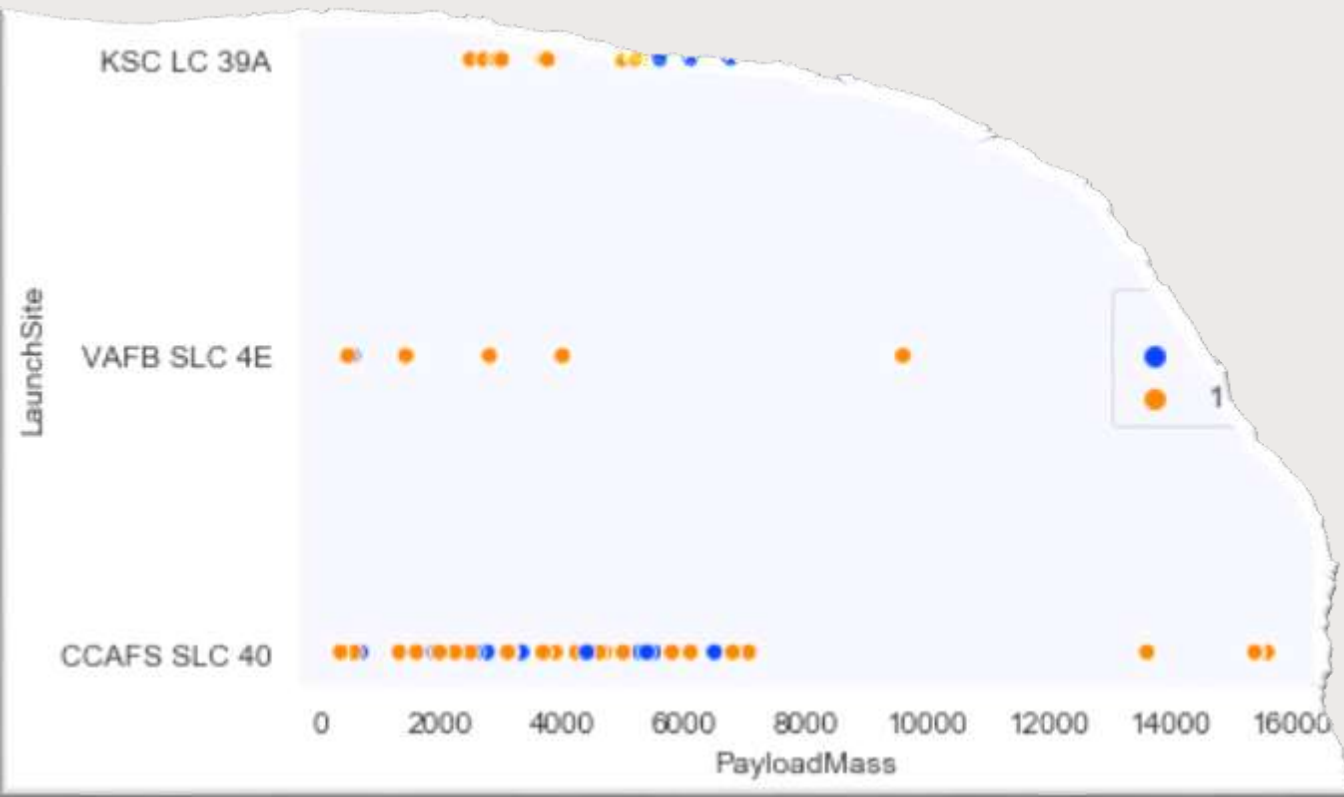
# **EDA with Visualization**

# Flight Number vs. Launch Site

- The scatter plot of Launch Site vs. Flight Number shows that:
- As the number of flights increases, the rate of success at a launch site increases.
- Most of the early flights (flight numbers < 30) were launched from CCAFS SLC 40, and were generally unsuccessful.
- The flights from VAFB SLC 4E also show this trend, that earlier flights were less successful.
- No early flights were launched from KSC LC 39A, so the launches from this site are more successful.
- Above a flight number of around 30, there are significantly more successful landings (Class = 1).



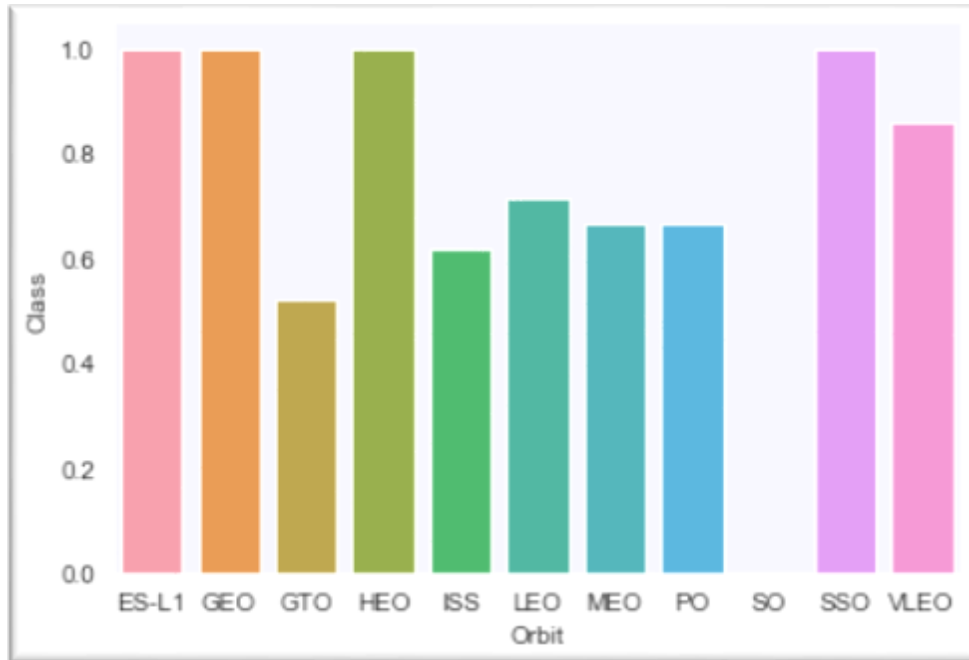
# Payload vs. Launch Site



- The screenshot of the scatter plot with explanation The scatter plot of Launch Site vs. Payload Mass shows that:
  - Above a payload mass of around 7000 kg, there are very few unsuccessful landings, but there is also far less data for these heavier launches.
  - There is no clear correlation between payload mass and success rate for a given launch site.
  - All sites launched a variety of payload masses, with most of the launches from CCAFS SLC 40 being comparatively lighter payloads (with some outliers).



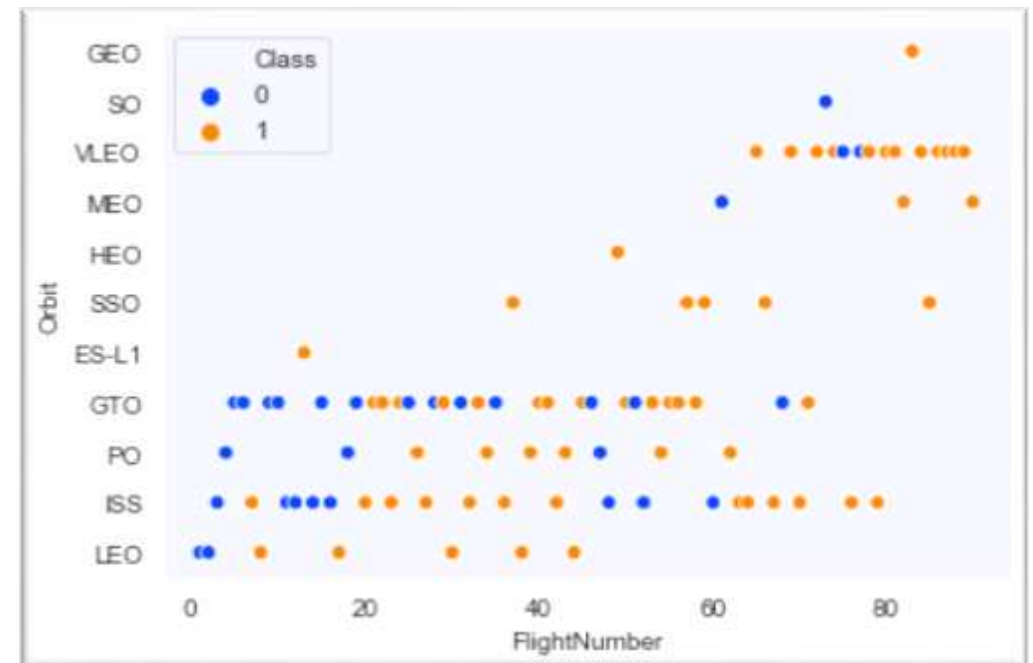
# Success Rate vs. Orbit Type



- The bar chart of Success Rate vs. Orbit Type shows that the following orbits have the highest (100%) success rate:
  - ES-L1 (Earth-Sun First Lagrangian Point)
  - GEO (Geostationary Orbit)
  - HEO (High Earth Orbit)
  - SSO (Sun-synchronous Orbit)
- The orbit with the lowest (0%) success rate is:
  - SO (Heliocentric Orbit)

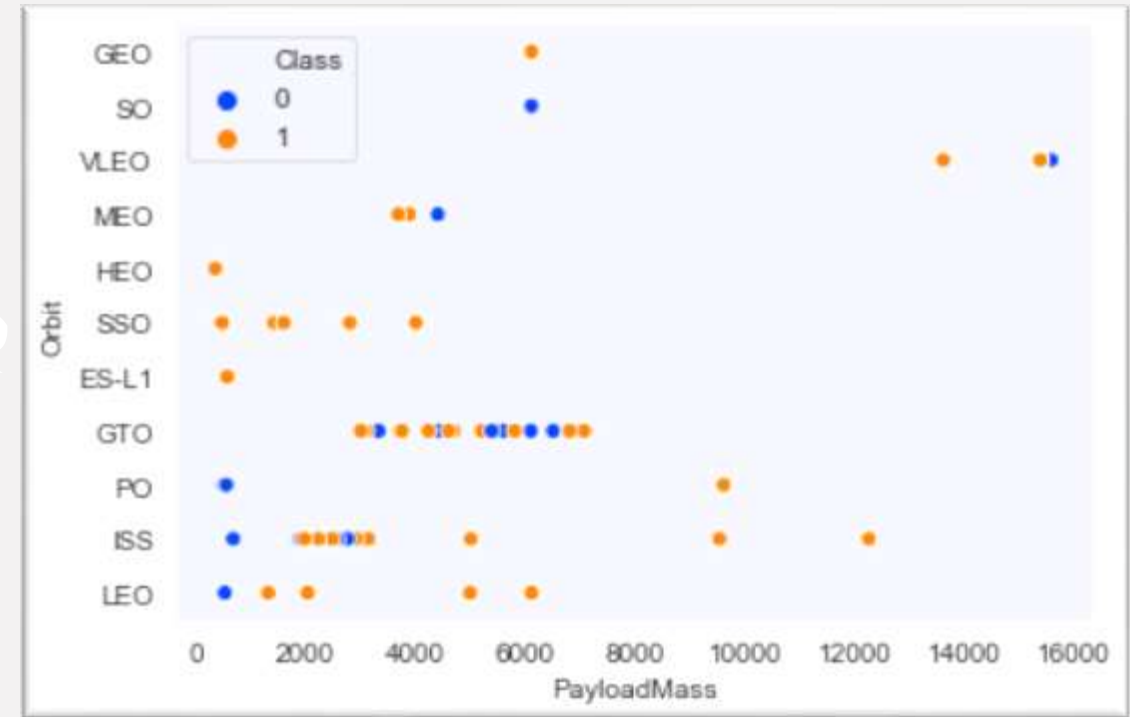
# Flight Number vs. Orbit Type

- This scatter plot of Orbit Type vs. Flight number shows a few useful things that the previous plots did not, such as:
  - The 100% success rate of GEO, HEO, and ES-L1 orbits can be explained by only having 1 flight into the respective orbits.
  - The 100% success rate in SSO is more impressive, with 5 successful flights.
  - There is little relationship between Flight Number and Success Rate for GTO.
  - Generally, as Flight Number increases, the success rate increases. This is most extreme for LEO, where unsuccessful landings only occurred for the low flight numbers (early launches).



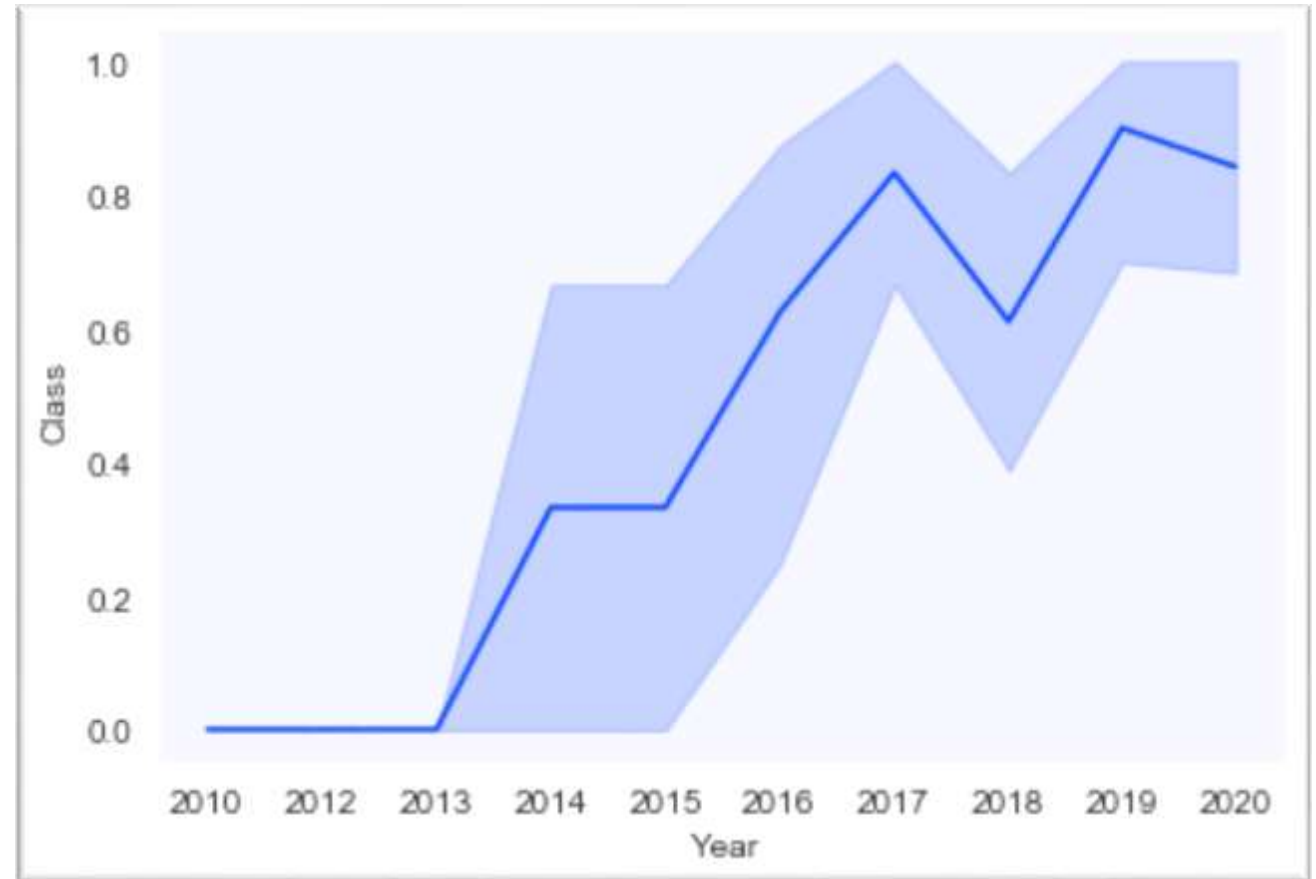
# Pay Load Mass vs. Orbit Type

- This scatter plot of Orbit Type vs. Payload Mass shows that:
  - The following orbit types have more success with heavy payloads:
    - PO (although the number of data points is small)
    - ISS
    - LEO
  - For GTO, the relationship between payload mass and success rate is unclear.
  - VLEO (Very Low Earth Orbit) launches are associated with heavier payloads, which makes intuitive sense.



# Launch Success Yearly Trend

- The line chart of yearly average success rate shows that:
  - Between 2010 and 2013, all landings were unsuccessful (as the success rate is 0).
  - After 2013, the success rate generally increased, despite small dips in 2018 and 2020.
  - After 2016, there was always a greater than 50% chance of success.



# EDA with SQL



# All Launch Site Names

- Find the names of the unique launch sites.

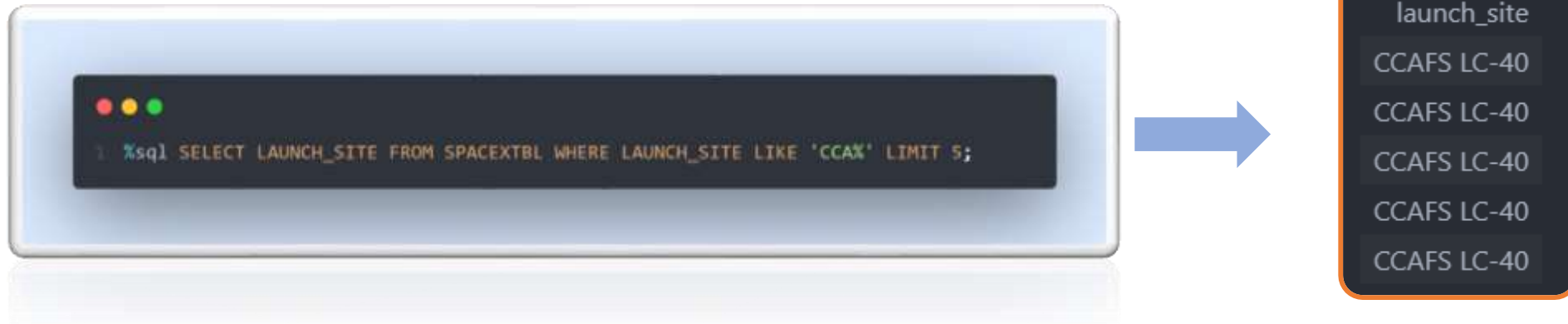


launch_site
CCAFS LC-40
CCAFS SLC-40
KSC LC-39A
VAFB SLC-4E

- The word `UNIQUE` returns only unique values from the `LAUNCH_SITE` column of the `SPACEXTBL` table.

# Launch Site Names Begin with 'CCA'

- Find 5 records where launch sites begin with 'CCA'.



- `LIMIT 5` fetches only 5 records, and the `LIKE` keyword is used with the wild card `'CCA%'` to retrieve string values beginning with 'CCA'.

# Total Payload Mass

- Calculate the total payload carried by boosters from NASA.

```
1 %sql SELECT SUM(PAYLOAD_MASS_KG_) AS TOTAL_PAYLOAD_MASS FROM SPACEXTBL \
2 WHERE CUSTOMER = 'NASA (CRS)';
```

total\_payload\_mass  
45596

- The **SUM** keyword is used to calculate the total of the **LAUNCH** column, and the **SUM** keyword (and the associated condition) filters the results to only boosters from NASA (CRS).

# Average Payload Mass by F9 v1.1

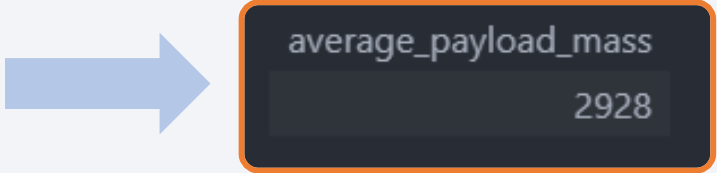
---

- Calculate the average payload mass carried by booster version F9 v1.1.



```
1 %sql SELECT AVG(PAYLOAD_MASS__KG_) AS AVERAGE_PAYLOAD_MASS FROM SPACEXTBL \
2 WHERE BOOSTER_VERSION = 'F9 v1.1';
```

The image shows a terminal window with a dark background and three colored window control buttons (red, yellow, green) in the top left corner. The SQL query is displayed in a light blue font. The terminal is part of a larger light blue rounded rectangle with a subtle reflection below it.



```
average_payload_mass
2928
```

A blue arrow points from the terminal window to a dark grey rounded rectangle with an orange border. Inside this rectangle, the result of the query is shown in two lines: the column name 'average\_payload\_mass' and the value '2928'.

- The **AVG** keyword is used to calculate the average of the **PAYLOAD\_MASS\_\_KG\_** column, and the **WHERE** keyword (and the associated condition) filters the results to only the F9 v1.1 booster version.

# First Successful Ground Landing Date

---

- Find the dates of the first successful landing outcome on ground pad.

```
1 %sql SELECT MIN(DATE) AS FIRST_SUCCESSFUL_GROUND_LANDING FROM SPACEXTBL \
2 WHERE LANDING__OUTCOME = 'Success (ground pad)';
```

first\_successful\_ground\_landing  
2015-12-22

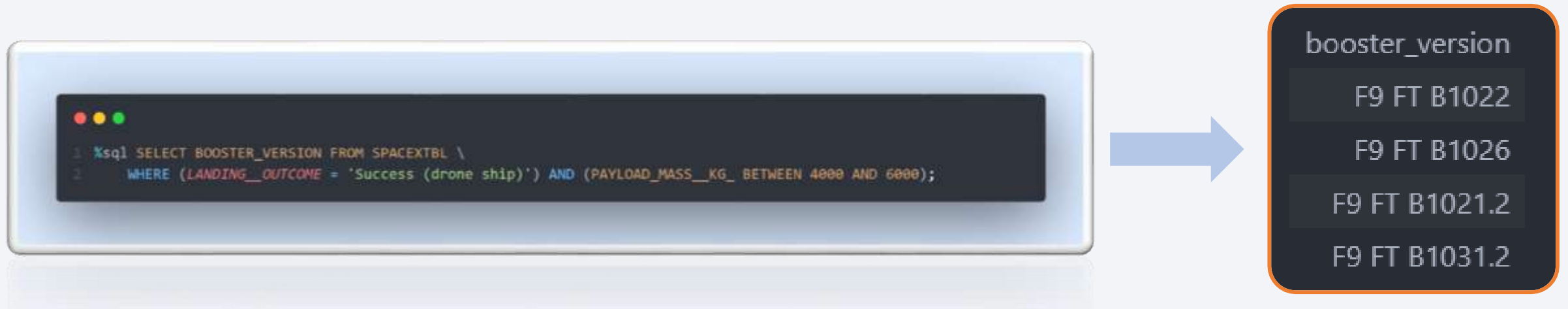
- The **MIN** keyword is used to calculate the minimum of the **DATE** column, i.e. the first date, and the **WHERE** keyword (and the associated condition) filters the results to only the successful ground pad landings.



# Successful Drone Ship Landing with Payload between 4000 and 6000

---

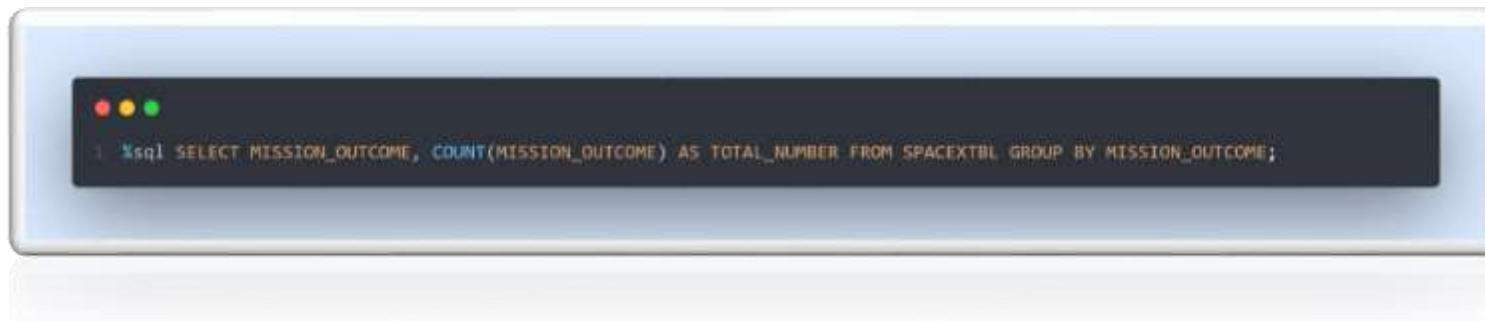
- List the names of boosters which have successfully landed on drone ship and had payload mass greater than 4000 but less than 6000.



- The **WHERE** keyword is used to filter the results to include only those that satisfy both conditions in the brackets (as the **AND** keyword is also used). The **BETWEEN** keyword allows for  $4000 < x < 6000$  values to be selected.

# Total Number of Successful and Failure Mission Outcomes

- Calculate the total number of successful and failure mission outcome.

A terminal window with a light blue border and a dark background. It contains a SQL query: `%sql SELECT MISSION_OUTCOME, COUNT(MISSION_OUTCOME) AS TOTAL_NUMBER FROM SPACEXTBL GROUP BY MISSION_OUTCOME;`. The query is preceded by a prompt character `%`.

```
%sql SELECT MISSION_OUTCOME, COUNT(MISSION_OUTCOME) AS TOTAL_NUMBER FROM SPACEXTBL GROUP BY MISSION_OUTCOME;
```



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

- The **COUNT** keyword is used to calculate the total number of mission outcomes, and the **GROUPBY** keyword is also used to group these results by the type of mission outcome.

# Boosters Carried Maximum Payload

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

```
1 %sql SELECT DISTINCT(BOOSTER_VERSION) FROM SPACEXTBL \
2 WHERE PAYLOAD_MASS_KG_ = (SELECT MAX(PAYLOAD_MASS_KG_) FROM SPACEXTBL);
```



booster\_version

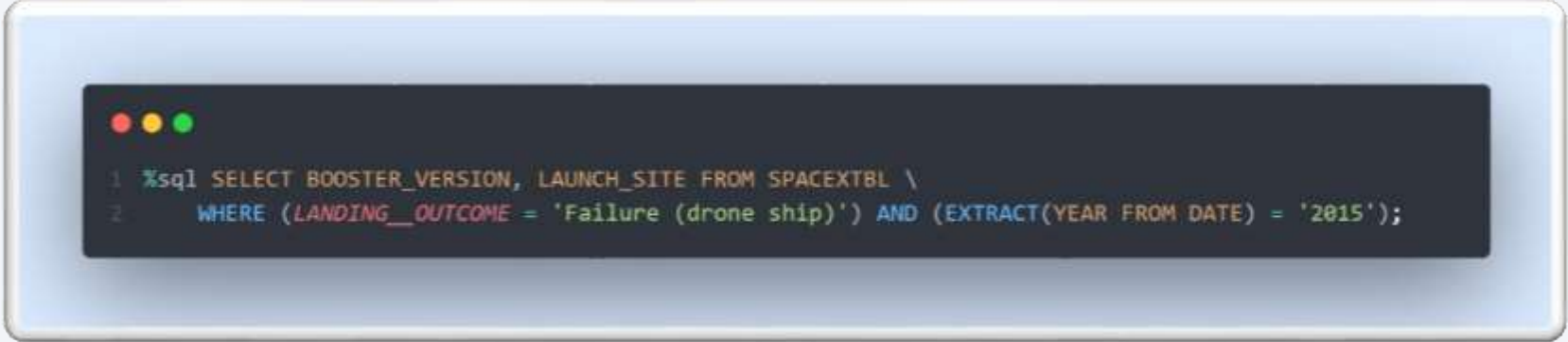
- F9 B5 B1048.4
- F9 B5 B1048.5
- F9 B5 B1049.4
- F9 B5 B1049.5
- F9 B5 B1049.7
- F9 B5 B1051.3
- F9 B5 B1051.4
- F9 B5 B1051.6
- F9 B5 B1056.4
- F9 B5 B1058.3
- F9 B5 B1060.2
- F9 B5 B1060.3

- A subquery is used here. The **SELECT** statement within the brackets finds the maximum payload mass and this value is used in the **WHERE** condition. The **DISTINCT** keyword is then used to retrieve only distinct /unique booster versions.


# 2015 Launch Records

---

- List the failed landing outcomes in drone ship, their booster versions, and launch site names for in year 2015.

A terminal window with a dark background and three colored window control buttons (red, yellow, green) in the top left corner. It contains two lines of SQL code.

```
1 %sql SELECT BOOSTER_VERSION, LAUNCH_SITE FROM SPACEXTBL \
2 WHERE (LANDING_OUTCOME = 'Failure (drone ship)') AND (EXTRACT(YEAR FROM DATE) = '2015');
```

A thick orange arrow pointing from the terminal window to the result table.

booster_version	launch_site
F9 v1.1 B1012	CCAFS LC-40
F9 v1.1 B1015	CCAFS LC-40

- The **WHERE** keyword is used to filter the results for only failed landing outcomes, **AND** only for the year of 2015.

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

- 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.

```
1 %sql SELECT LANDING_OUTCOME, COUNT(LANDING_OUTCOME) AS TOTAL_NUMBER FROM SPACEXTBL \
2 WHERE DATE BETWEEN '2010-06-04' AND '2017-03-20' \
3 GROUP BY LANDING_OUTCOME \
4 ORDER BY TOTAL_NUMBER DESC;
```



landing_outcome	total_number
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 **WHERE** keyword is used with the **BETWEEN** keyword to filter the results to dates only within those specified. The results are then grouped and ordered, using the keywords **GROUP BY** and **ORDER BY**, respectively, where **DESC** is used to specify the descending order.

**LAUNCH SITES**

**PROXIMITY ANALYSIS**

**– FOLIUM**

**INTERACTIVE MAP**

# LAUNCH SITES PROXIMITY ANALYSIS – FOLIUM INTERACTIVE MAP

---



All SpaceX launch sites are on coasts of the United States of America, specifically Florida and California.



# SUCCESS/FAILED LAUNCHES FOR EACH SITE

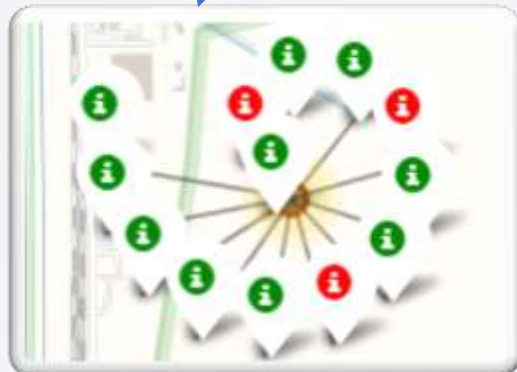
---



VAFB SLC-4E



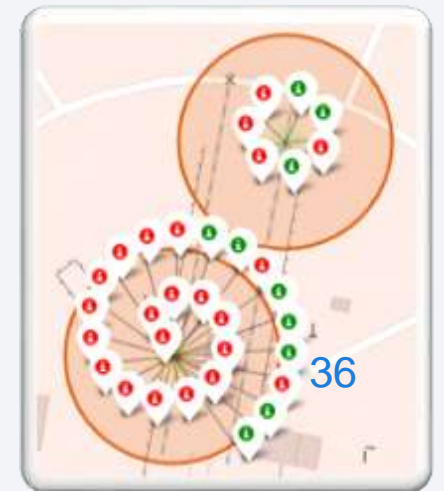
KSC LC-39A



CCAFS SLC-40 and CCAFS LC-40

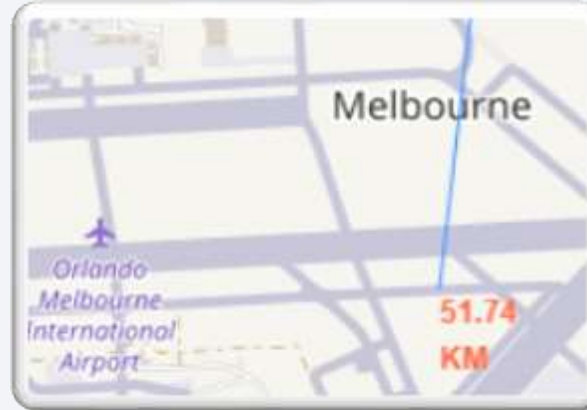


=



# PROXIMITY OF LAUNCH SITES TO OTHER POINTS OF INTEREST

Using the CCAFS SLC-40 launch site as an example site, we can understand more about the placement of launch sites.



Are launch sites in close proximity to railways?

- YES. The coastline is only 0.87 km due East.

Are launch sites in close proximity to highways?

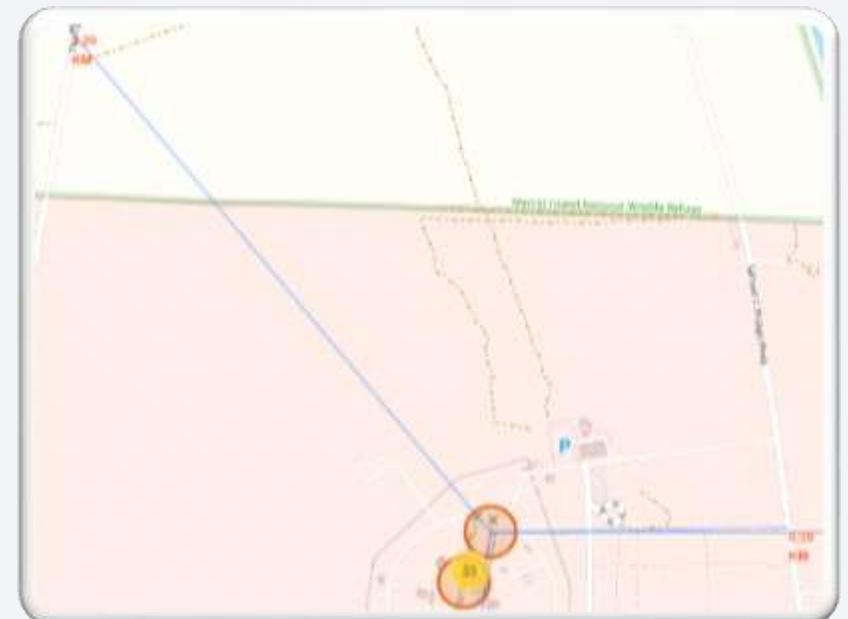
- YES. The nearest highway is only 0.59km away.

Are launch sites in close proximity to railways?

- YES. The nearest railway is only 1.29 km away.

Do launch sites keep certain distance away from cities?

- YES. The nearest city is 51.74 km away.

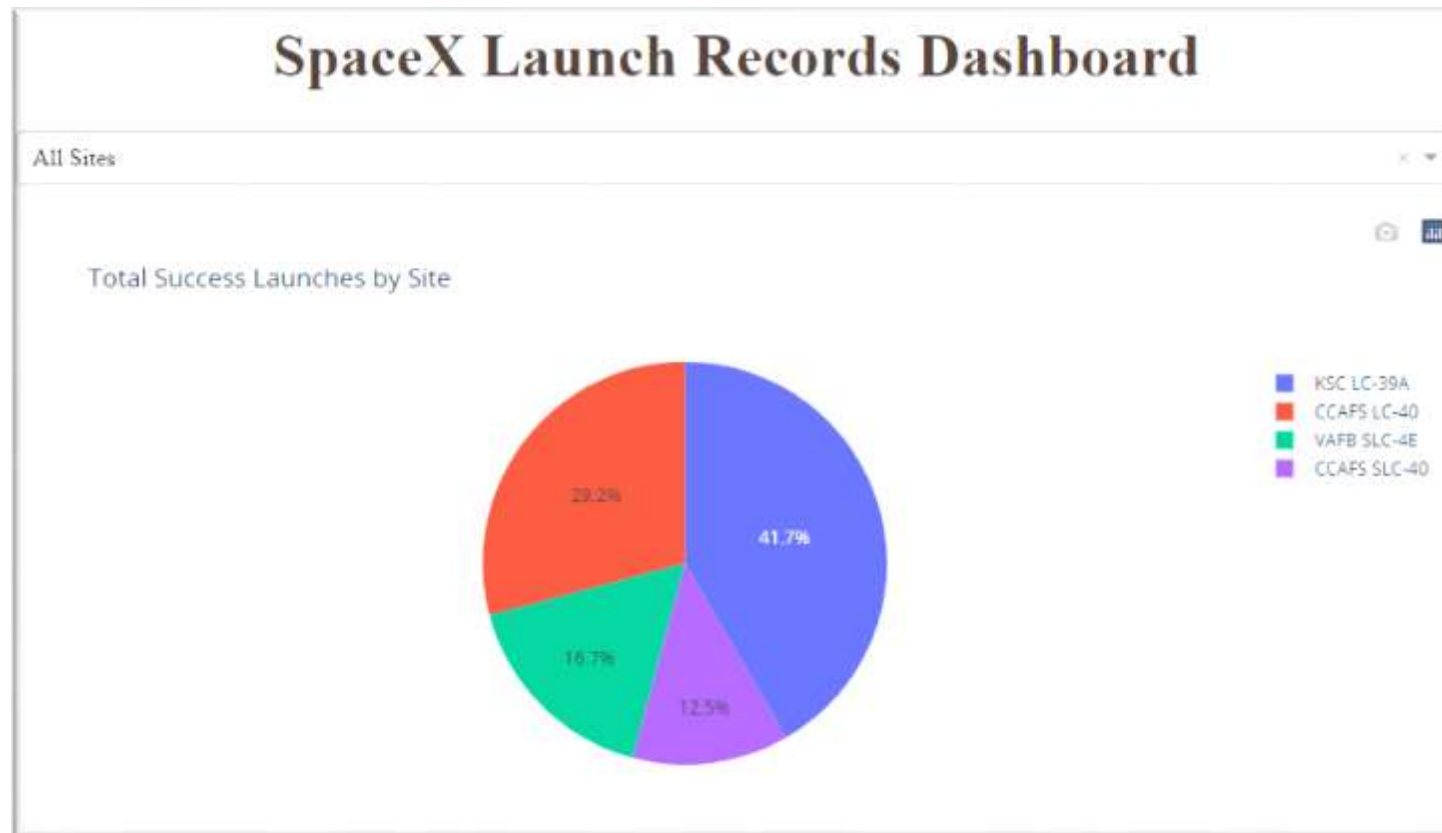


# interactive dashboard

## - Plotly Dash

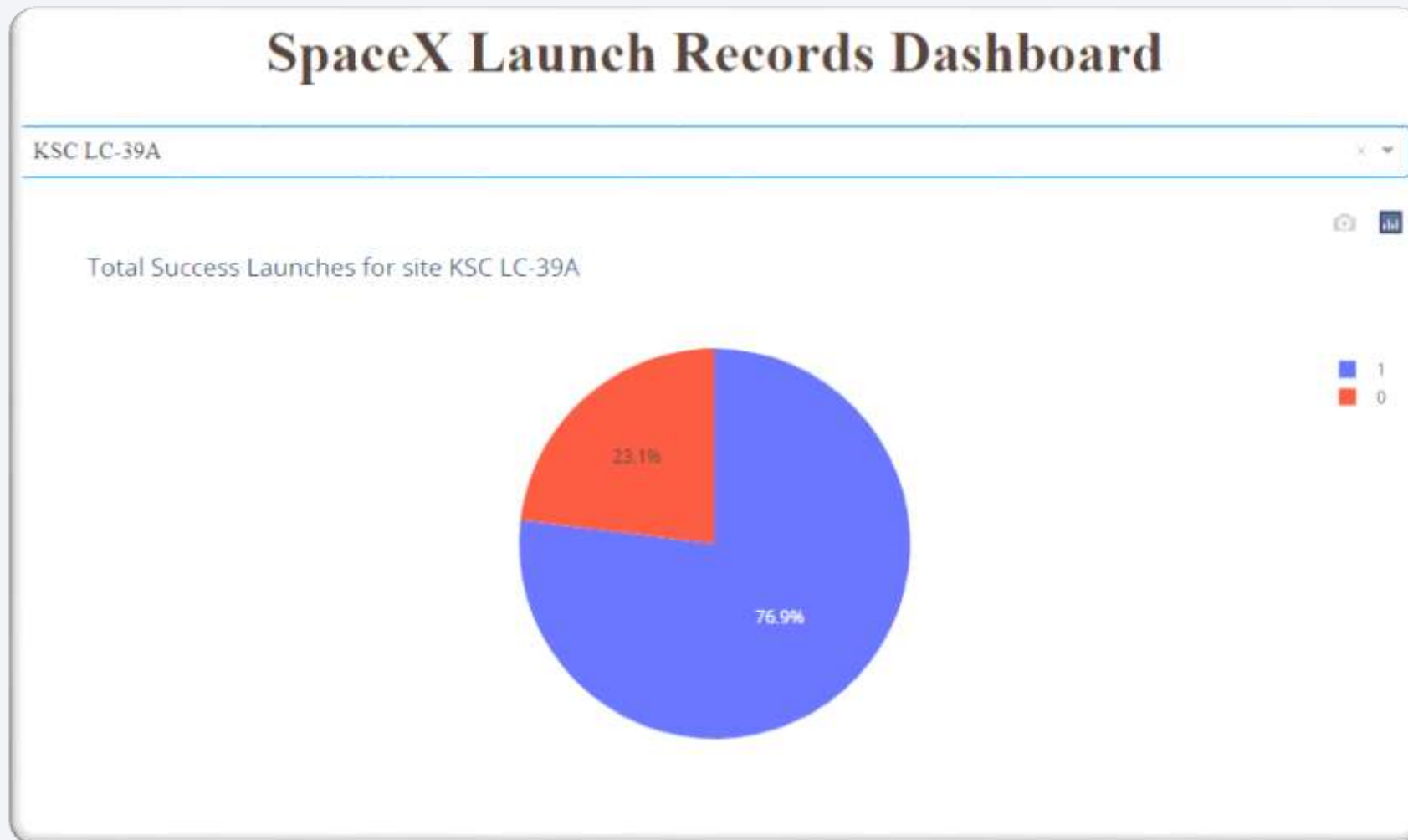
launch success  
count for all sites

- The launch site KSC LC-39 A had the most successful launches, with 41.7% of the total successful launches.



## Pie chart for the launch site with highest launch success ratio

---



The launch site KSC LC-39 A also had the highest rate of successful launches, with a 76.9% success rate.

# Launch Outcome vs. Payload scatter plot for all sites



- Plotting the launch outcome vs. payload for all sites shows a gap around 4000 kg, so it makes sense to split the data into 2 ranges:
  - 0 – 4000 kg (low payloads)
  - 4000 – 10000 kg (massive payloads)
- From these 2 plots, it can be shown that the success for massive payloads is lower than that for low payloads.
- It is also worth noting that some booster types (v1.0 and B5) have not been launched with massive payloads.



# **PREDICTIVE ANALYSIS**

---

## **- CLASSIFICATION**

---

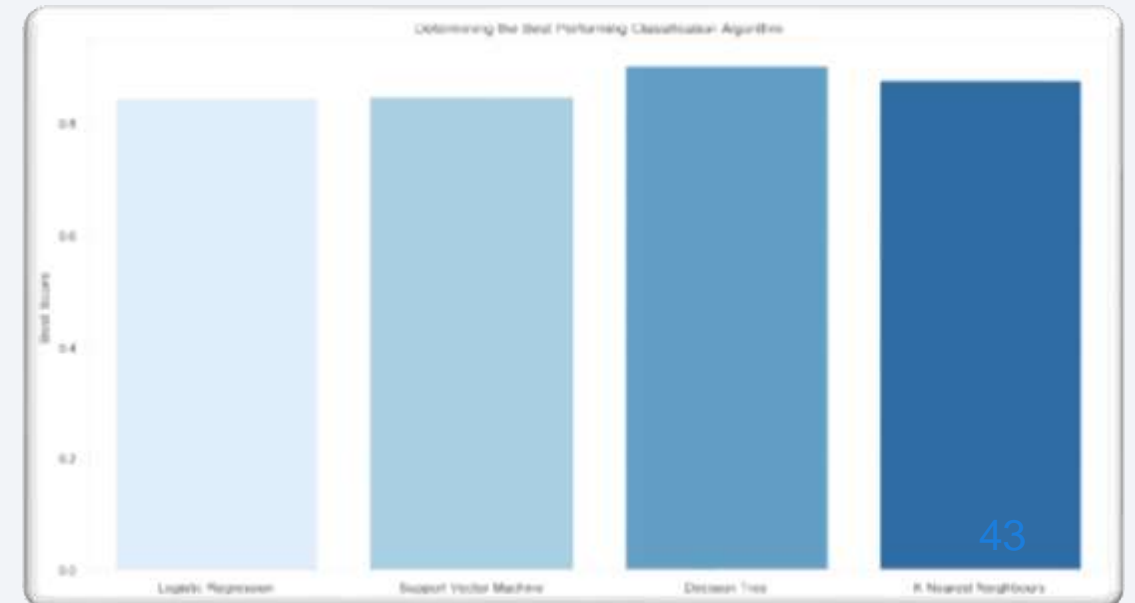
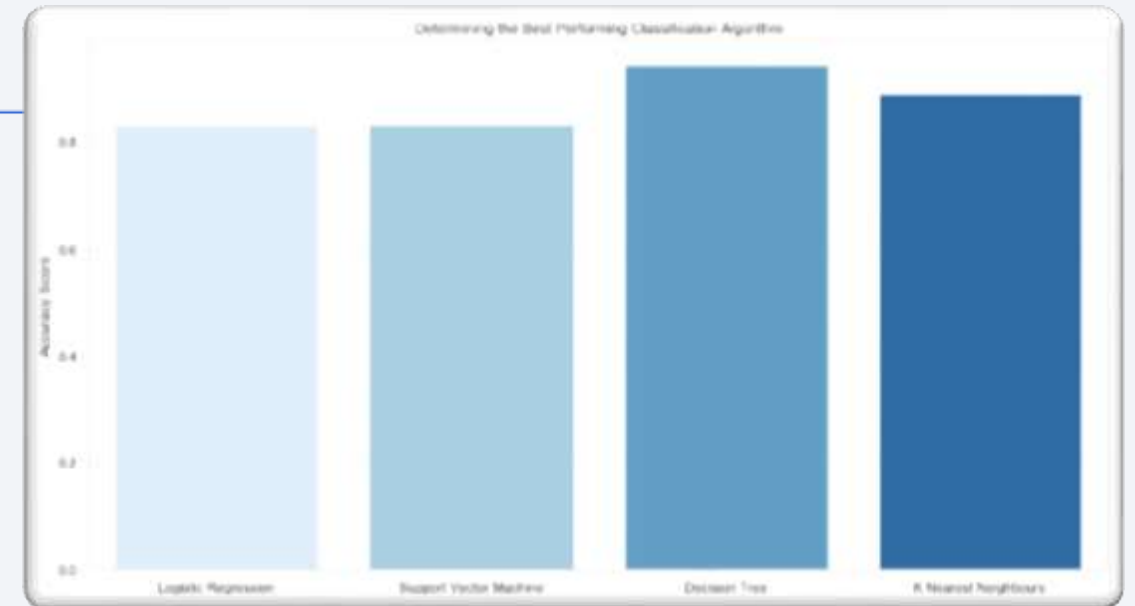


# Classification Accuracy

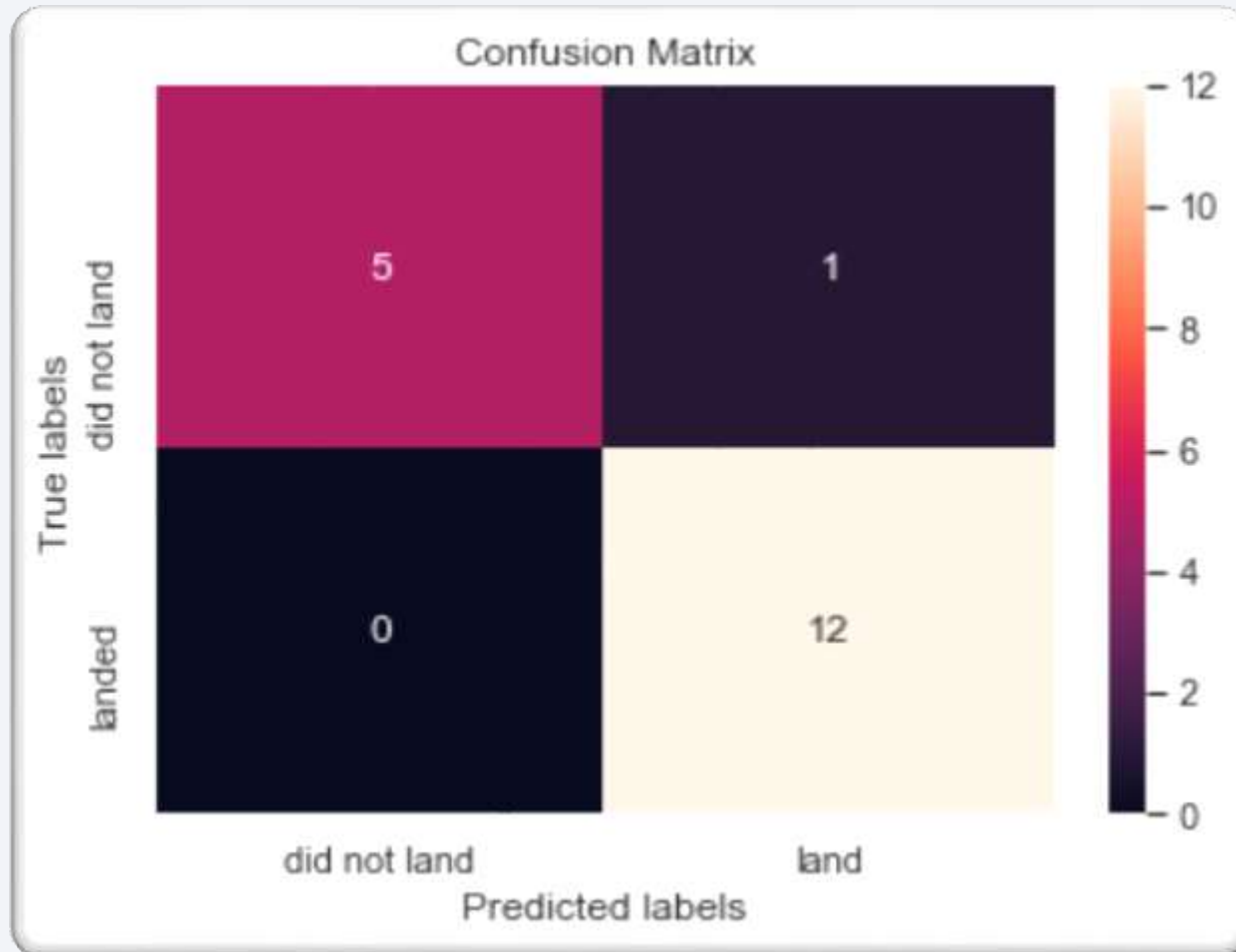
Plotting the Accuracy Score and Best Score for each classification algorithm produces the following result:

- The **Decision Tree** model has the highest classification accuracy
- The Accuracy Score is 94.44%
- The Best Score is 90.36%

Algorithm	Accuracy Score	Best Score
Logistic Regression	0.833333	0.846429
Support Vector Machine	0.833333	0.848214
Decision Tree	0.944444	0.903571
K Nearest Neighbours	0.888889	0.876786



# Confusion Matrix



- As shown previously, best performing classification model is the **Decision Tree** model, with an accuracy of 94.44%.
- This is explained by the confusion matrix, which shows only 1 out of 18 total results classified incorrectly (a false positive, shown in the top-right corner).
- The other 17 results are correctly classified (5 did not land, 12 did land).

# CONCLUSIONS

# CONCLUSIONS

- As the number of flights increases, the rate of success at a launch site increases, with most early flights being unsuccessful. I.e. with more experience, the success rate increases.
  - Between 2010 and 2013, all landings were unsuccessful (as the success rate is 0).
  - After 2013, the success rate generally increased, despite small dips in 2018 and 2020.
  - After 2016, there was always a greater than 50% chance of success.
- Orbit types ES-L1, GEO, HEO, and SSO, have the highest (100%) success rate.
  - The 100% success rate of GEO, HEO, and ES-L1 orbits can be explained by only having 1 flight into the respective orbits.
  - The 100% success rate in SSO is more impressive, with 5 successful flights.
  - The orbit types PO, ISS, and LEO, have more success with heavy payloads:
  - VLEO (Very Low Earth Orbit) launches are associated with heavier payloads, which makes intuitive sense.
- The launch site KSC LC-39 A had the most successful launches, with 41.7% of the total successful launches, and also the highest rate of successful launches, with a 76.9% success rate.
- The success for massive payloads (over 4000kg) is lower than that for low payloads.
- The best performing classification model is the Decision Tree model, with an accuracy of 94.44%.



# Appendix

## DATA COLLECTION – space x REST api

```
From the rocket column we would like to learn the booster name.

# Takes the dataset and uses the rocket column to call the API and append the data to the list
def getBoosterVersion(data):
    for x in data['rocket']:
        response = requests.get("https://api.spacexdata.com/v4/rockets/"+str(x)).json()
        BoosterVersion.append(response['name'])

from the launchpad we would like to know the name of the launch site being used, the longitude, and the latitude.

# Takes the dataset and uses the launchpad column to call the API and append the data to the list
def getLaunchSite(data):
    for x in data['launchpad']:
        response = requests.get("https://api.spacexdata.com/v4/launchpads/"+str(x)).json()
        Longitude.append(response['longitude'])
        Latitude.append(response['latitude'])
        LaunchSite.append(response['name'])

From the payload we would like to learn the mass of the payload and the orbit that it is going to.

# Takes the dataset and uses the payload column to call the API and append the data to the list
def getPayloadData(data):
    for load in data['payloads']:
        response = requests.get("https://api.spacexdata.com/v4/payloads/"+load).json()
        PayloadMass.append(response['mass_kg'])
        Orbit.append(response['orbit'])
```

```
From cores we would like to learn the outcome of the landing, the type of the landing, number of flights with that core, whether gridfins were used, whether the core is reused, whether legs were used, the landing pad used, the block of the core (which is a number used to separate versions of cores), the number of times this specific core has been reused, and the serial of the core.

# Takes the dataset and uses the cores column to call the API and append the data to the list
def getCoresData(data):
    for core in data['cores']:
        if core['core'] != None:
            response = requests.get("https://api.spacexdata.com/v4/cores/"+core['core']).json()
            Block.append(response['block'])
            ReusedCount.append(response['reuse_count'])
            Serial.append(response['serial'])
        else:
            Block.append(None)
            ReusedCount.append(None)
            Serial.append(None)
        Outcome.append(str(core['landing_success'])+" "+str(core['landing_type']))
        Flights.append(core['flights'])
        GridFins.append(core['gridfins'])
        Reused.append(core['reused'])
        Legs.append(core['legs'])
        LandingPad.append(core['landing_pad'])
```

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

# We will remove rows with multiple cores because there are Falcon rockets with 2 active reusert boosters
# and some that have multiple payloads in a single rocket
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 remove the single values in the list and replace the features
data['cores'] = data['cores'].map(lambda x : x[0])
data['payloads'] = data['payloads'].map(lambda x : x[0])

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

# Using the date we will restrict the dates of the launches
data = data[data['date'] >= datetime.date(2010, 10, 10)]
```

- Custom functions to retrieve the required information
- Custom logic to clean the data

# Appendix

## DATA COLLECTION – WEB SCRAPING

- Custom functions for web scraping
- Custom logic to fill up the launch\_dict values with values from the launch tables

```
def data_time(table_cells):
    """
    This function returns the date and time from the HTML table cell
    Input: the element of a table data cell extracts extra row
    """
    return [data_time.strip() for data_time in list(table_cells.strings)][0:2]

def booster_version(table_cells):
    """
    This function returns the booster version from the HTML table cell
    Input: the element of a table data cell extracts extra row
    """
    out = ''
    for i, booster_version in enumerate(table_cells.strings):
        if i % 2 == 0:
            out = out + booster_version
    return out

def landing_status(table_cells):
    """
    This function returns the landing status from the HTML table cell
    Input: the element of a table data cell extracts extra row
    """
    out = [1 for i in table_cells.strings]
    return out

def get_mass(table_cells):
    mass = unicode(data.normalize("WWT", table_cells.text).strip())
    if mass:
        mass.find("kg")
        new_mass = mass[0:mass.find("kg")+2]
    else:
        new_mass = 0
    return new_mass

def extract_column_from_header(row):
    """
    This function returns the landing status from the HTML table cell
    Input: the element of a table data cell extracts extra row
    """
    if (row.br):
        row.br.extract()
    if row.a:
        row.a.extract()
    if row.sup:
        row.sup.extract()

    column_name = ' '.join(row.contents)

    # Filter the digit and empty names
    if not(column_name.strip().isdigit()):
        column_name = column_name.strip()
        return column_name
```

```
extracted_row = {}
# Extract each table
for table_number, table in enumerate(soup.find_all("table", {"class": "table table-bordered table-striped"})):
    # Get table row
    for row in table.find_all("tr"):
        # Detect to see if first table heading to be a number corresponding to launch a number
        if row.th:
            if row.th.string:
                flight_number = row.th.string.strip()
                flight_number.isdigit()
            else:
                # Flag=False
                # Get table element
                row = row.find_all("tr")
                if 0 in launch_row.index:
                    # Flag
                    extracted_row = {}
                    # Flight Number value
                    # Append the flight number into launch_dict with key 'Flight No.'
                    launch_dict["Flight No."] = append(flight_number)

                    # Date value
                    # Append the date into launch_dict with key 'Date'
                    date = row[1].string
                    data = date.split(" ")
                    launch_dict["Date"] = append(date)

                    # Time value
                    # Append the time into launch_dict with key 'Time'
                    time = row[2].string
                    launch_dict["Time"] = append(time)

                    # Booster version
                    # Append the booster version into launch_dict with key 'Booster version'
                    booster_version = row[3].string
                    if booster_version:
                        launch_dict["Booster version"] = append(booster_version)

                    # Launch site
                    # Append the launch site into launch_dict with key 'Launch site'
                    launch_site = row[4].string
                    launch_dict["Launch site"] = append(launch_site)

                    # Payload
                    # Append the payload into launch_dict with key 'Payload'
                    payload = row[5].string
                    launch_dict["Payload"] = append(payload)

                    # Payload mass
                    # Append the payload mass into launch_dict with key 'Payload mass'
                    payload_mass = row[6].string
                    launch_dict["Payload mass"] = append(payload_mass)

                    # Price
                    # Append the price into launch_dict with key 'Price'
                    price = row[7].string
                    launch_dict["Price"] = append(price)

                    # Customer
                    # Append the customer into launch_dict with key 'Customer'
                    if row[8].string:
                        customer = row[8].string
                    else:
                        customer = "None"
                    launch_dict["Customer"] = append(customer)

                    # Launch outcome
                    # Append the launch outcome into launch_dict with key 'Launch outcome'
                    launch_outcome = row[9].string
                    launch_dict["Launch outcome"] = append(launch_outcome)

                    # Booster landing
                    # Append the booster landing into launch_dict with key 'Booster landing'
                    booster_landing = row[10].string
                    launch_dict["Booster landing"] = append(booster_landing)

    print("Flight Number: ", flight_number, " Date: ", date, " Time: ", time, " \n",
          "Booster Version: ", booster_version, " Launch Site: ", launch_site, " \n",
          "Payload: ", payload, " Price: ", price, " \n",
          "Customer: ", customer, " Launch Outcome: ", launch_outcome, " \n",
          "Booster Landing: ", booster_landing, " \n",
          " \n")
```

Thank you!

