

Winning Space Race with Data Science

Lawrence Chinn
18 Nov 2024



Outline

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

Executive Summary

Summary of methodologies

1. Data collection methodology
 - Request from SpaceX API and web scraping from Wikipedia of Falcon Launches (https://en.wikipedia.org/wiki/List_of_Falcon_9_and_Falcon_Heavy_launches)
2. Perform data wrangling
 - The two data sets were cleaned, transformed and integrated into a single set for quality control before EDA
3. Perform exploratory data analysis (EDA) using visualization and SQL
4. Perform interactive visual analytics using Folium and Plotly Dash
5. Perform predictive analysis using classification models
 - Build and assess multiple training models

Summary of all results

1. SpaceX developed a degree of expertise over time (2010-2020) allowed them to achieve a relative high level of reliability when launching their Falcon F9 rockets
2. However, the Falcon F9 is limited by its payload restriction and some orbits (e.g. GTO) remain a challenge.

Introduction

Background

In the span of 10 years, SpaceX's accomplishments include

- Sending a spacecraft to the International Space Station
- Establishing Starlink, a satellite-based network providing Internet access
- Manned missions to Space.

For SpaceX, the rocket launches are relatively inexpensive. SpaceX advertises Falcon 9 rocket launches on its website for \$62 million, whereas other providers advertise cost upwards of \$165 million. Much of the savings is because, unlike other providers, SpaceX can reuse the first stage. However, this is not always the case. Sometimes, it will crash and is sacrificed due to the mission parameters, e.g. payload, orbit, and customer.

Problem

Accurately determine the successful recovery of the Falcon 9 booster to gain some insight for the actual cost of a rocket launch.

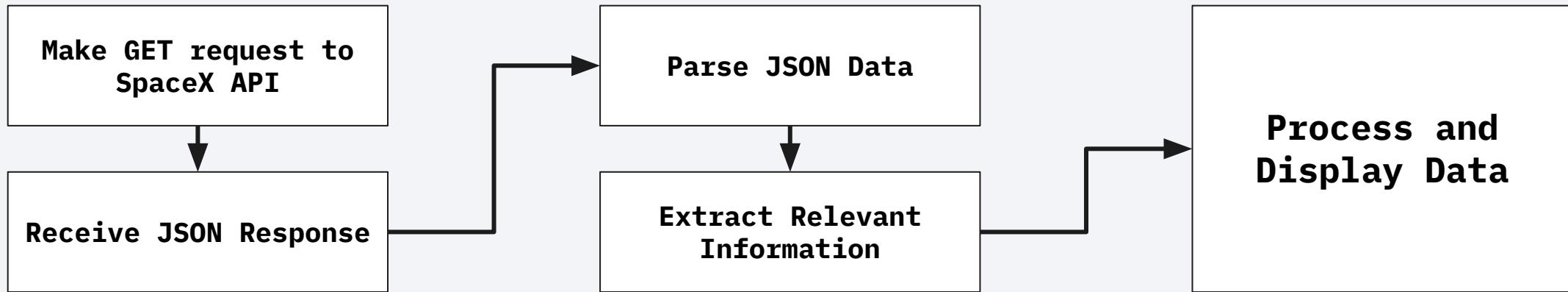
Section 1

Methodology

Methodology

- Data collection methodology
 - Request from SpaceX API and web scraping from Wikipedia of Falcon Launches (https://en.wikipedia.org/wiki/List_of_Falcon_9_and_Falcon_Heavy_launches)
- Perform data wrangling
 - The two data sets were cleaned, transformed and integrated into a single set for quality control before EDA
- Perform exploratory data analysis (EDA) using visualization and SQL
- Perform interactive visual analytics using Folium and Plotly Dash
- Perform predictive analysis using classification models
 - Build and assess multiple training models

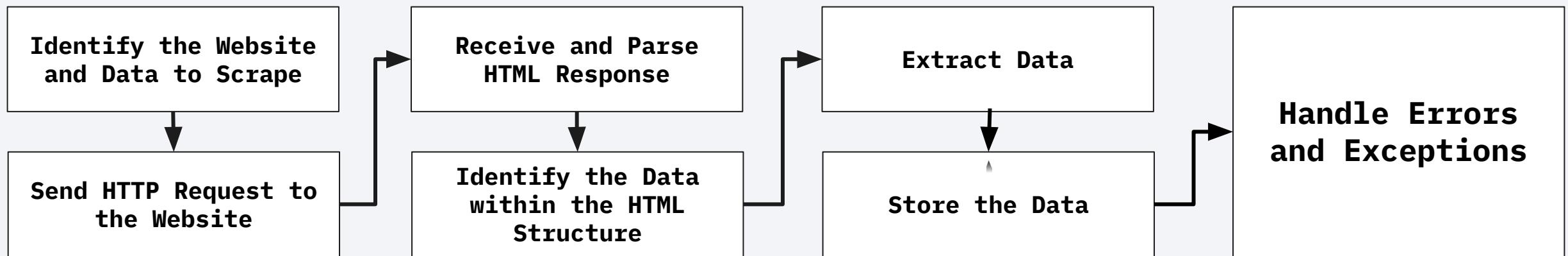
Data Collection - SpaceX API



GitHub URL:

https://github.com/blarghwtfbbq/Capstone-_SpaceX/blob/c454b1b52469f4981f2c4ec047159cff5944be75/jupyter-labs-spacex-data-collection-api%20Mod%201A.ipynb

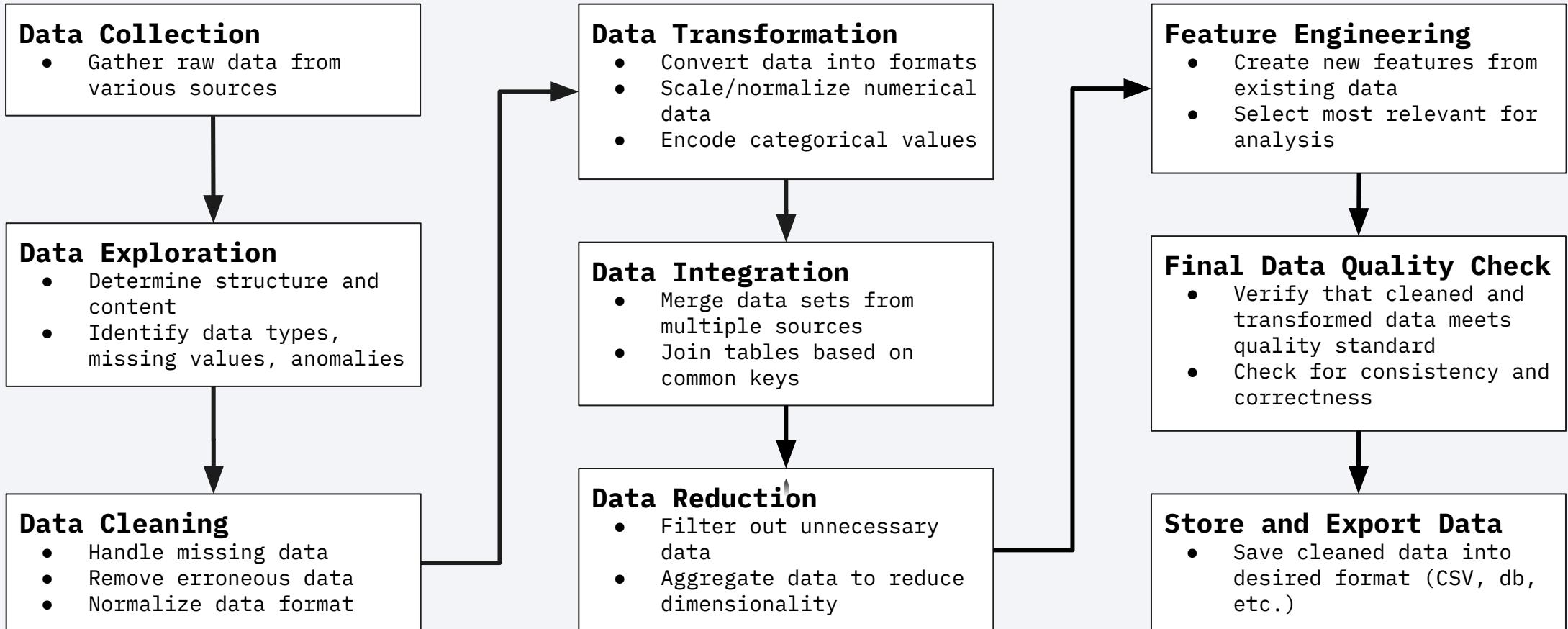
Data Collection - Scraping



GitHub URL:

https://github.com/blarughwtfbbq/Capstone-_SpaceX/blob/c454b1b52469f4981f2c4ec047159cff5944be75/jupyter-labs-webscraping%20Mod%201B.ipynb

Data Wrangling



9

GitHub URL:

https://github.com/blarghwtfbq/Capstone-_SpaceX/blob/c454b1b52469f4981f2c4ec047159cff5944be75/labs-jupyter-spacex-Data%20wrangling%20Mod%201C.ipynb

EDA with Data Visualization

Visualizations

- Scatter plot: Flight number vs Launch Site
- Scatter plot: Payload Mass vs Launch Site
- Bar chart: Success rate by Orbit Type
- Scatter plot: Flight Number vs Orbit Type
- Scatter plot: Payload Type vs Orbit Type
- Line plot: Yearly Launch Success Rate Trend

GitHub URL:

https://github.com/blarughwtfbbq/Capstone_-_SpaceX/blob/c454b1b52469f4981f2c4ec047159cff5944be75/edadataviz%20Mod%20B.ipynb

EDA with SQL

SQL Queries

- Names of unique launch sites
- Five 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
- Date when first successful landing outcome in ground pad was achieved
- List names of boosters which have success in drone ship and have payload greater than 4000 kg but less than 6000 kg
- List total number of successful and failure mission outcomes
- List names of the booster version which carried the maximum payload mass (subquery)
- List records which display month names, failure outcomes in drone ship, booster versions, launch site for the months in the year 2015
- Rank the count landing outcomes between 2010-06-04 and 2017-03-20 in descending order

GitHub URL:

https://github.com/blarughwtfbbq/Capstone-_SpaceX/blob/c454b1b52469f4981f2c4ec047159cff5944be75/jupyter-labs-eda-sql-coursera_sqlite%20Mod%202A.ipynb

Build an Interactive Map with Folium

Map Features

1. Markers to identify locations of SpaceX launch sites across the U.S.
2. Sub-markers for the number of launches that are color-coded to indicate success or failure
3. Distance calculations to indicate proximity to support infrastructure and populated areas

GitHub URL:

https://github.com/blarughwtfbbq/Capstone-_SpaceX/blob/c454b1b52469f4981f2c4ec047159cff5944be75/lab_jupyter_launch_site_location%20Mod%203A.ipynb

Build a Dashboard with Plotly Dash

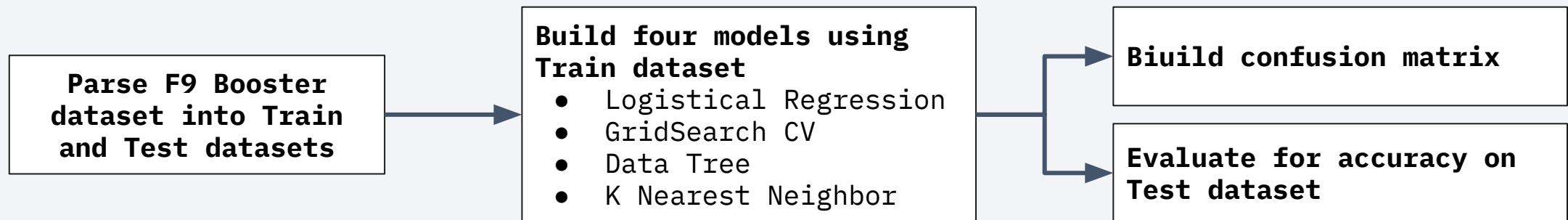
Dashboard Features

1. Range slider: Allows selection of payload mass that range from 0 kg up to 10,000 kg by increments of 2500 kg
2. Pie Chart callback: Updates based upon selected launch site
3. Scatter Chart callback: Updates scatter chart based upon selected launch site and payload range
 - a. Correlate launch success and payload mass to color-labeled dots representing booster version category

GitHub URL:

https://github.com/blarughwtfbbq/Capstone-_SpaceX/blob/c454b1b52469f4981f2c4ec047159cff5944be75/spacex_dash_app.py

Predictive Analysis (Classification)

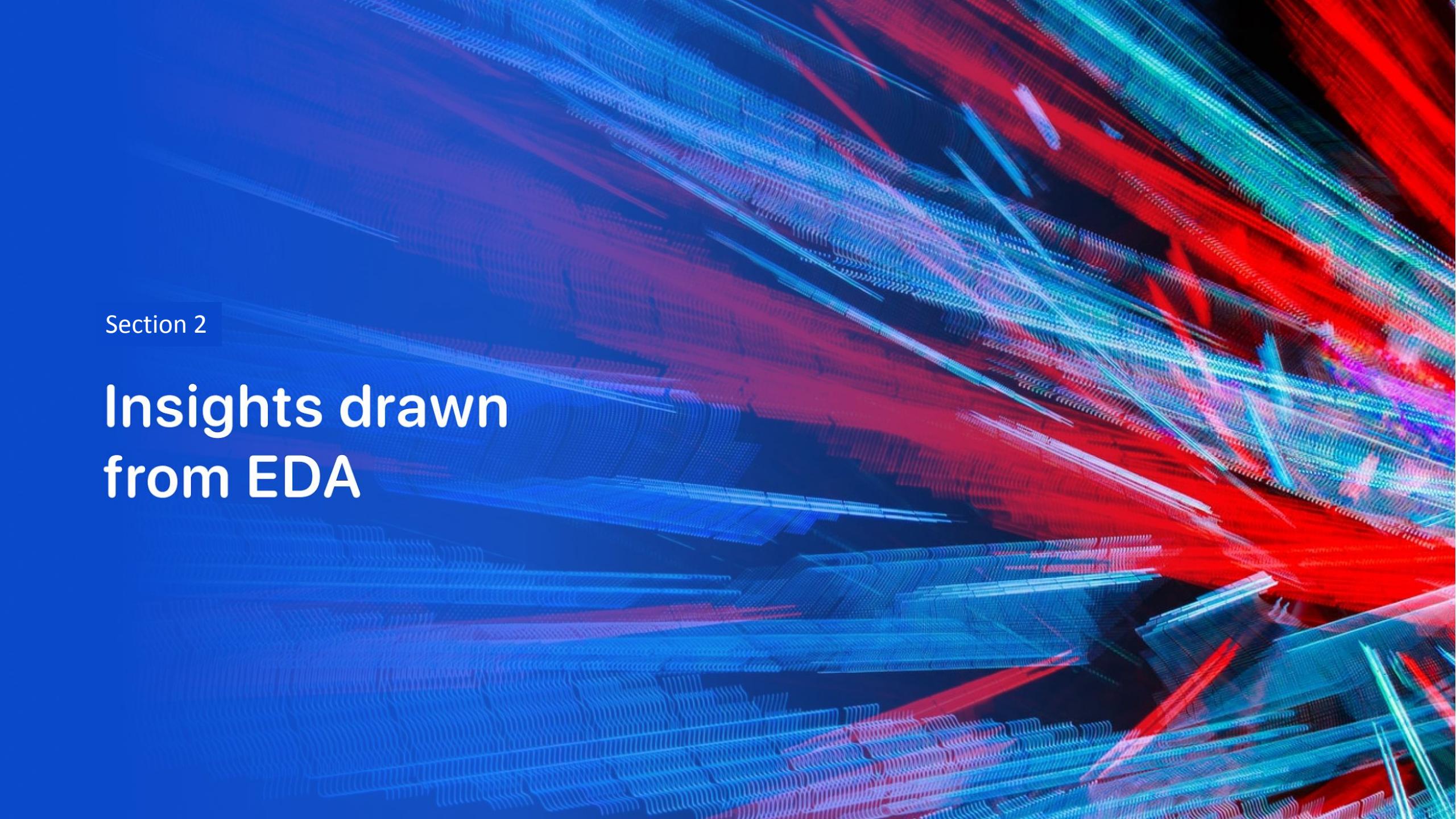


GitHub URL:

https://github.com/blarughwtfbbq/Capstone_SpaceX/blob/c56ee2e1c810bb89c5e34fcc983d683616332543/SpaceX_Machine_Learning_Prediction_Part_5_v1.ipynb

Results

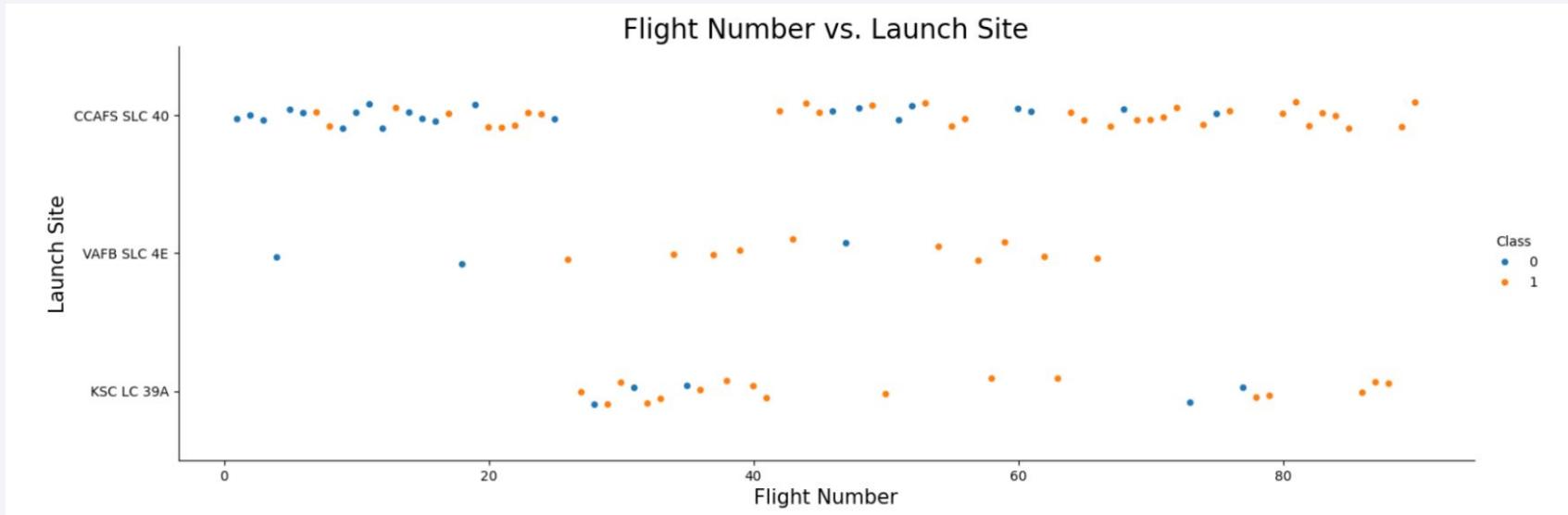
- Exploratory results suggest that SpaceX developed and maintain a high degree of expertise in an iterative fashion allowing them to maintain a relatively stable level of successful launches between 2010 and 2020.
- Interactive analytics demo in screenshots
- Predictive analysis results

The background of the slide features a complex, abstract pattern of glowing lines. These lines are primarily blue and red, creating a sense of depth and motion. They appear to be composed of numerous small, glowing particles or dots, giving them a textured, almost organic appearance. The lines converge and diverge, forming various shapes and directions across the dark, solid-colored background.

Section 2

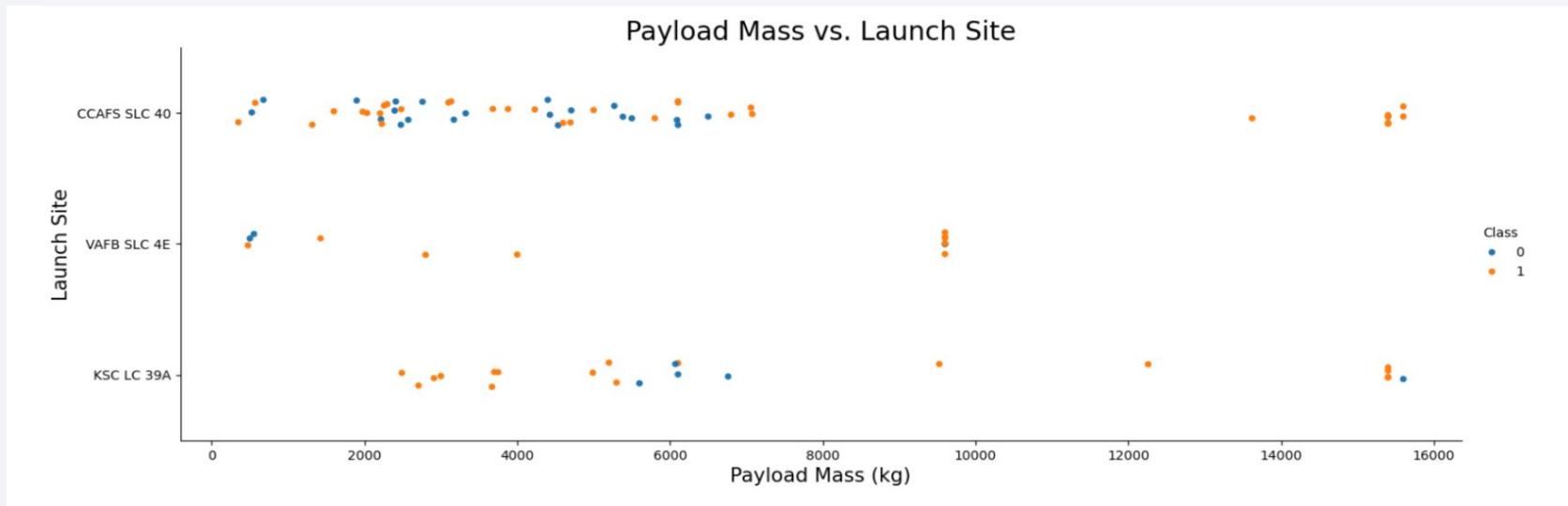
Insights drawn from EDA

Flight Number vs. Launch Site



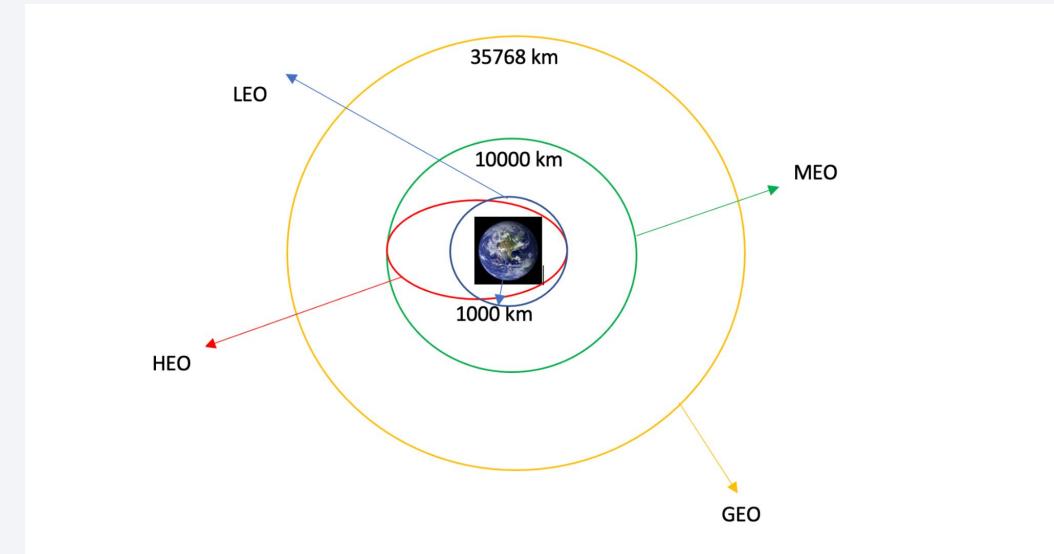
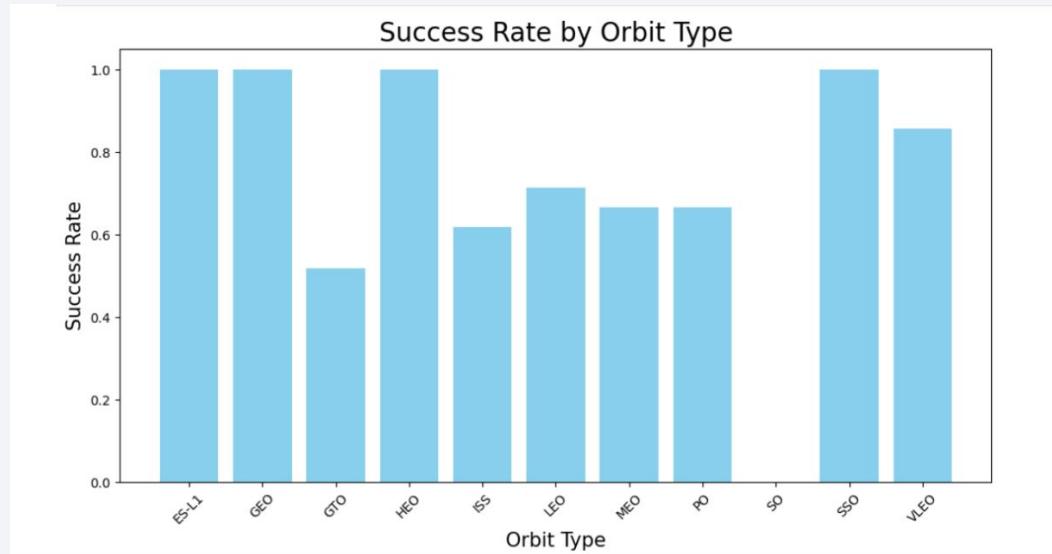
1. CCAFS SLC-40 may appear more frequently indicating it may be designated as a primary location for launches
2. Shift from failure to success with increasing flight number may be attributed improvement in technology, experience, or strategies
3. Greater success observed at CCAFS SLC-40 and KSC LC-39A may be due to factors such as site conditions, technology compatibility, or mission-specific complexities

Payload vs. Launch Site



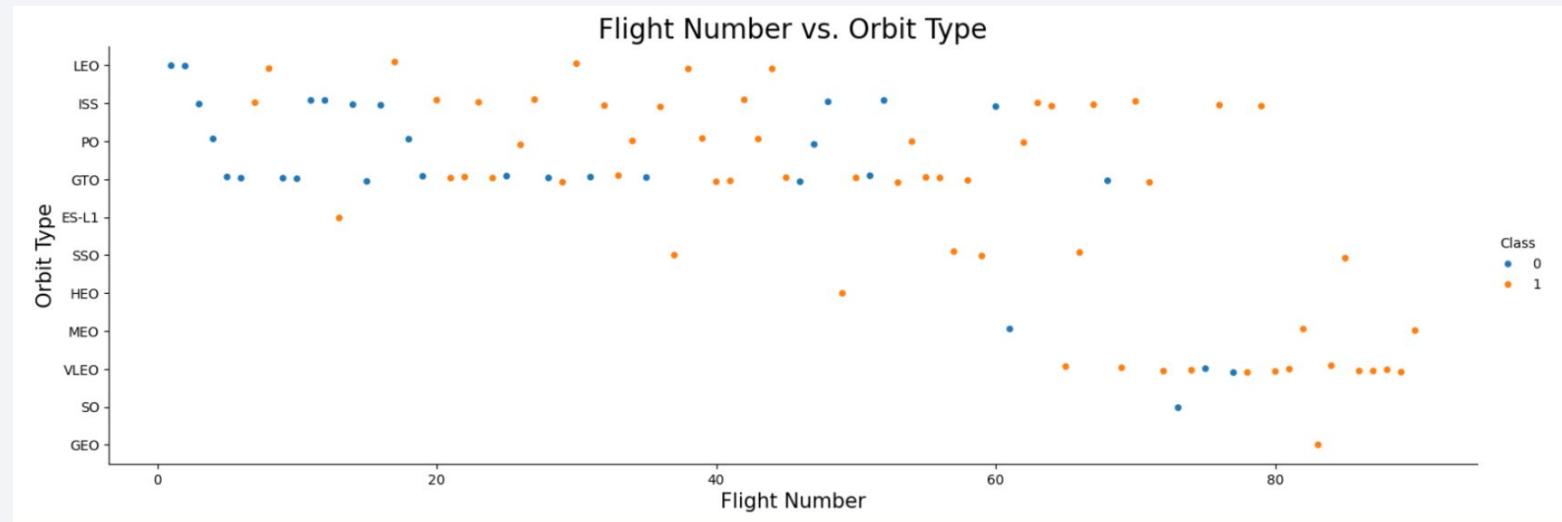
1. CCAFS SLC-40 may appear more frequently indicating it may be designated as a primary location for launches
2. Success pattern
 - o No rockets launched from VAFB-SLC launch site with a heavy payload (greater than 10,000 kg)
 - o Certain site may exhibit consistent success for specific payload ranges due to optimal conditions or better technology
3. Key infrastructure or procedural strengths at those locations may play a role in successful launches with heavy payloads at CCAFS SLC-40 and KSC LC-39A

Success Rate vs. Orbit Type



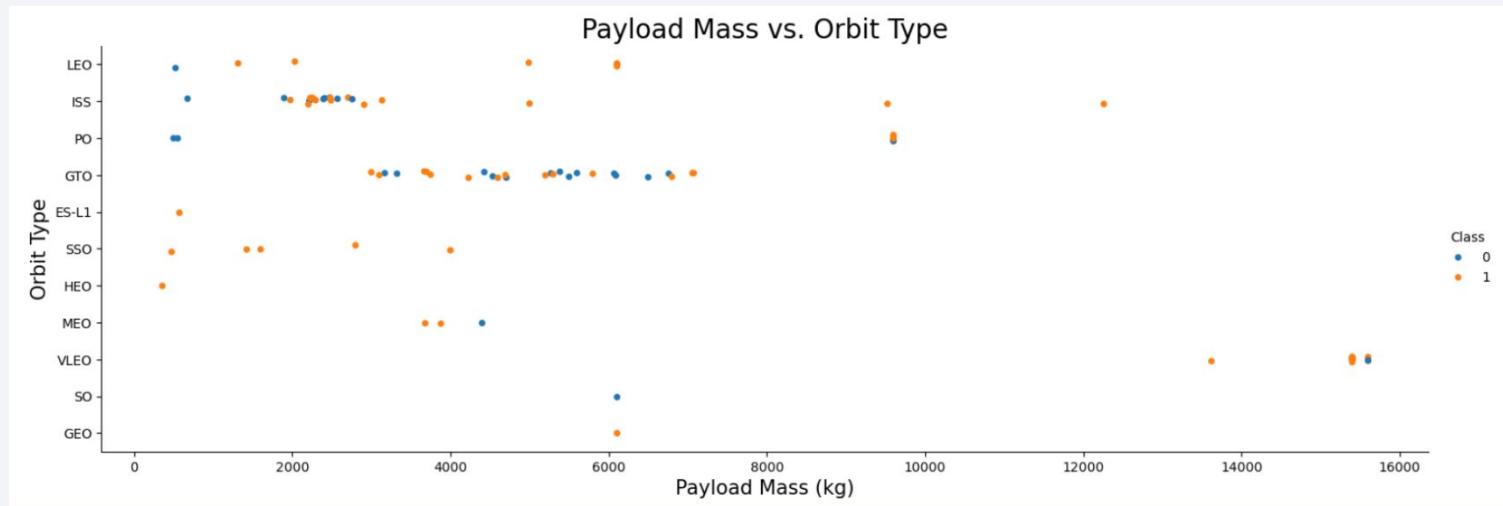
1. Higher success rates observed for launch attempts to the GEO, HEO, and VLEO
2. Lower success rates observed for launch attempts to GTO, LEO, and MEO
3. Generally, there seems to be a higher success rate for higher orbit attempts greater than 10,000 km
 - Further investigation is required to identify the challenges to successfully achieve LEO to MEO range as well as the failure to target GTO

Flight Number vs. Orbit Type



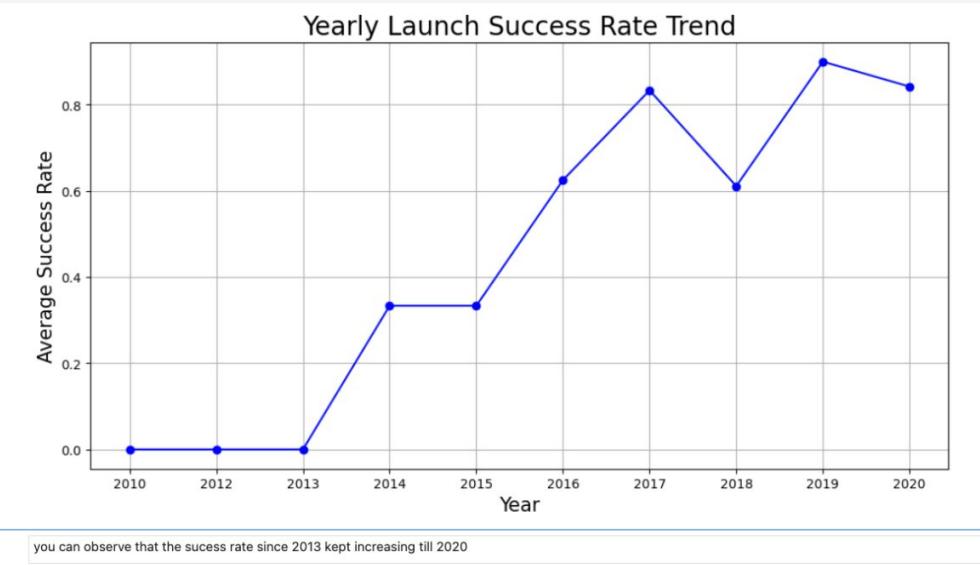
1. Success rate for LEO may improve with higher flight numbers, e.g. flight experience
2. Other factors, not related to flight numbers, have an impact upon the success for GTO

Payload vs. Orbit Type



1. Success was observed for LEO, ISS, and PO for heavy payloads
2. For GTO, the outcome is a mixture of success and failure across a range of payload sizes confounding any conclusion

Launch Success Yearly Trend



1. Positive trend in success rate with some yearly fluctuations possibly due to challenges, e.g. technical issues, adverse weather conditions
2. Pronounced improvement reflects its strive to continually learn and enhance reliability leading to a stabilized high level of success

SpaceX made significant improvement with its launch success rate over time. This can be attributed to a continuous effort to learn, innovate, and improve their launch systems

All Launch Site Names

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

```
# Write the SQL query to get unique launch sites
query = 'SELECT DISTINCT "Launch_Site" FROM SPACEXTBL'

# Execute the query and load the results into a DataFrame
unique_launch_sites = pd.read_sql_query(query, conn)

# Display the results
print(unique_launch_sites)
```

	Launch_Site
0	CCAFS LC-40
1	VAFB SLC-4E
2	KSC LC-39A
3	CCAFS SLC-40

```
Date Time (UTC) Booster_Version Launch_Site \
0 2010-06-04 18:45:00 F9 v1.0 B0003 CCAFS LC-40
1 2010-12-08 15:43:00 F9 v1.0 B0004 CCAFS LC-40
2 2012-05-22 7:44:00 F9 v1.0 B0005 CCAFS LC-40
3 2012-10-08 0:35:00 F9 v1.0 B0006 CCAFS LC-40
4 2013-03-01 15:10:00 F9 v1.0 B0007 CCAFS LC-40

Payload PAYLOAD_MASS_KG_ \
0 Dragon Spacecraft Qualification Unit 0
1 Dragon demo flight C1, two CubeSats, barrel of... 0
2 Dragon demo flight C2 525
3 SpaceX CRS-1 500
4 SpaceX CRS-2 677

Orbit Customer Mission_Outcome Landing_Outcome
0 LEO SpaceX Success Failure (parachute)
1 LEO (ISS) NASA (COTS) NRO Success Failure (parachute)
2 LEO (ISS) NASA (COTS) Success No attempt
3 LEO (ISS) NASA (CRS) Success No attempt
4 LEO (ISS) NASA (CRS) Success No attempt
Index(['Date', 'Time (UTC)', 'Booster_Version', 'Launch_Site', 'Payload',
       'PAYLOAD_MASS_KG_', 'Orbit', 'Customer', 'Mission_Outcome',
       'Landing_Outcome'],
      dtype='object')
```

Launch Site Names Begin with 'CCA'

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

```
# SQL query to get 5 records where launch sites begin with 'CCA'
query = """
SELECT *
FROM SPACEXTBL
WHERE "Launch_Site" LIKE 'CCA%'
LIMIT 5;
"""

# Execute the query and load the results into a DataFrame
cca_launch_sites = pd.read_sql_query(query, conn)

# Display the results
print(cca_launch_sites)
```

	Date	Time (UTC)	Booster_Version	Launch_Site	Payload	PAYLOAD_MASS__KG_
0	2010-06-04	18:45:00	F9 v1.0	B0003	CCAFS LC-40	0
1	2010-12-08	15:43:00	F9 v1.0	B0004	CCAFS LC-40	0
2	2012-05-22	7:44:00	F9 v1.0	B0005	CCAFS LC-40	525
3	2012-10-08	0:35:00	F9 v1.0	B0006	CCAFS LC-40	500
4	2013-03-01	15:10:00	F9 v1.0	B0007	CCAFS LC-40	677

	Orbit	Customer	Mission_Outcome	Landing_Outcome
0	LEO	SpaceX	Success	Failure (parachute)
1	LEO (ISS)	NASA (COTS) NRO	Success	Failure (parachute)
2	LEO (ISS)	NASA (COTS)	Success	No attempt
3	LEO (ISS)	NASA (CRS)	Success	No attempt
4	LEO (ISS)	NASA (CRS)	Success	No attempt

Total Payload Mass

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

```
# SQL query to get the total payload mass carried by boosters launched by NASA (CRS)
query = """
SELECT SUM("PAYLOAD_MASS__KG_") AS Total_Payload_Mass
FROM SPACEXTBL
WHERE "Customer" = 'NASA (CRS)';
"""

# Execute the query and load the results into a DataFrame
total_payload_mass = pd.read_sql_query(query, conn)

# Display the results
print(total_payload_mass)
```

	Total_Payload_Mass
0	45596

Average Payload Mass by F9 v1.1

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

```
# SQL query to get the average payload mass carried by booster version F9 v1.1
query = """
SELECT AVG("PAYLOAD_MASS__KG_") AS Average_Payload_Mass
FROM SPACEXTBL
WHERE "Booster_Version" = 'F9 v1.1';
"""

# Execute the query and load the results into a DataFrame
average_payload_mass = pd.read_sql_query(query, conn)

# Display the results
print(average_payload_mass)
```

	Average_Payload_Mass
0	2928.4

First Successful Ground Landing Date

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

Hint:Use min function

```
# SQL query to find the date of the first successful landing on a ground pad
query = """
SELECT MIN("Date") AS First_Successful_Landing_Date
FROM SPACEXTBL
WHERE "Landing_Outcome" = 'Success (ground pad)';
"""

# Execute the query and load the results into a DataFrame
first_successful_landing_date = pd.read_sql_query(query, conn)

# Display the results
print(first_successful_landing_date)
```

First_Successful_Landing_Date
0 2015-12-22

Successful Drone Ship Landing with Payload between 4000 and 6000

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

```
# SQL query to list booster versions with specific criteria
query = """
SELECT "Booster_Version"
FROM SPACEXTBL
WHERE "Landing_Outcome" = 'Success (drone ship)'
AND "PAYLOAD_MASS__KG_" > 4000
AND "PAYLOAD_MASS__KG_" < 6000;
"""

# Execute the query and load the results into a DataFrame
boosters_success_drone_ship = pd.read_sql_query(query, conn)

# Display the results
print(boosters_success_drone_ship)
```

	Booster_Version
0	F9 FT B1022
1	F9 FT B1026
2	F9 FT B1021.2
3	F9 FT B1031.2

Total Number of Successful and Failure Mission Outcomes

List the total number of successful and failure mission outcomes

```
# SQL query to get the total number of successful and failure mission outcomes
query = """
SELECT "Mission_Outcome", COUNT(*) AS Outcome_Count
FROM SPACEXTBL
GROUP BY "Mission_Outcome";
"""

# Execute the query and load the results into a DataFrame
mission_outcomes = pd.read_sql_query(query, conn)

# Display the results
print(mission_outcomes)
```

	Mission_Outcome	Outcome_Count
0	Failure (in flight)	1
1	Success	98
2	Success	1
3	Success (payload status unclear)	1

Boosters Carried Maximum Payload

List the names of the booster_versions which have carried the maximum payload mass. Use a subquery

```
# SQL query to list booster versions with the maximum payload mass using a subquery
query = """
SELECT "Booster_Version"
FROM SPACEXTBL
WHERE "PAYLOAD_MASS_KG_" = (
    SELECT MAX("PAYLOAD_MASS_KG_")
    FROM SPACEXTBL
);
"""

# Execute the query and load the results into a DataFrame
boosters_max_payload = pd.read_sql_query(query, conn)

# Display the results
print(boosters_max_payload)
```

	Booster_Version
0	F9 B5 B1048.4
1	F9 B5 B1049.4
2	F9 B5 B1051.3
3	F9 B5 B1056.4
4	F9 B5 B1048.5
5	F9 B5 B1051.4
6	F9 B5 B1049.5
7	F9 B5 B1060.2
8	F9 B5 B1058.3
9	F9 B5 B1051.6
10	F9 B5 B1060.3
11	F9 B5 B1049.7

2015 Launch Records

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, 6,2) as month to get the months and substr(Date,0,5)='2015' for year.

```
# SQL query to get records with specific criteria
query = '''
SELECT~
    SUBSTR("Date", 6, 2) AS Month,~
    "Landing_Outcome",~
    "Booster_Version",~
    "Launch_Site"
FROM SPACEXTBL
WHERE~
    "Landing_Outcome" = 'Failure (drone ship)'~
    AND SUBSTR("Date", 0, 5) = '2015';
'''

# Execute the query and load the results into a DataFrame
failure_landings_2015 = pd.read_sql_query(query, conn)

# Display the results
print(failure_landings_2015)
```

	Month	Landing_Outcome	Booster_Version	Launch_Site
0	01	Failure (drone ship)	F9 v1.1 B1012	CCAFS LC-40
1	04	Failure (drone ship)	F9 v1.1 B1015	CCAFS LC-40

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.

```
# SQL query to rank the count of landing outcomes between the specified dates in descending order
query = """
SELECT "Landing_Outcome", COUNT(*) AS Outcome_Count
FROM SPACEXTBL
WHERE "Date" BETWEEN '2010-06-04' AND '2017-03-20'
GROUP BY "Landing_Outcome"
ORDER BY Outcome_Count DESC;
"""

# Execute the query and load the results into a DataFrame
landing_outcomes_ranked = pd.read_sql_query(query, conn)

# Display the results
print(landing_outcomes_ranked)
```

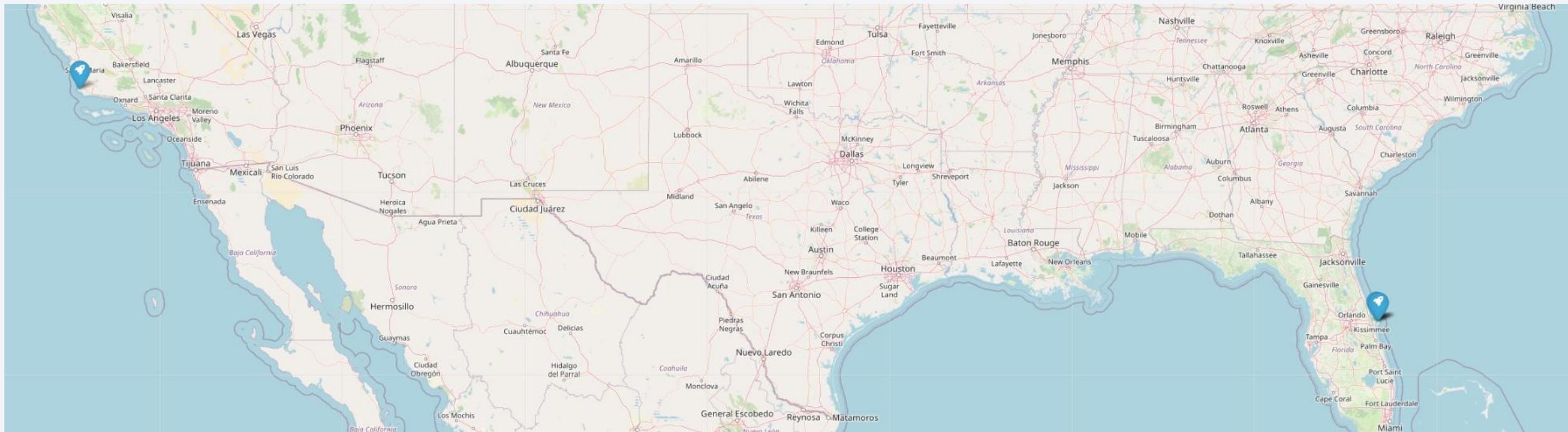
	Landing_Outcome	Outcome_Count
0	No attempt	10
1	Success (drone ship)	5
2	Failure (drone ship)	5
3	Success (ground pad)	3
4	Controlled (ocean)	3
5	Uncontrolled (ocean)	2
6	Failure (parachute)	2
7	Precluded (drone ship)	1

The background of the slide is a photograph taken from space at night. It shows the curvature of the Earth's horizon against a dark blue sky. City lights are visible as small white dots, with larger clusters of lights indicating major urban areas. In the upper right corner, there is a faint, greenish glow of the aurora borealis or a similar atmospheric phenomenon.

Section 3

Launch Sites Proximities Analysis

Location of SpaceX Launch Sites



1. Markers indicate that the SpaceX launch sites were located in California and Florida

Florida Launch sites

CCAFS LC-40

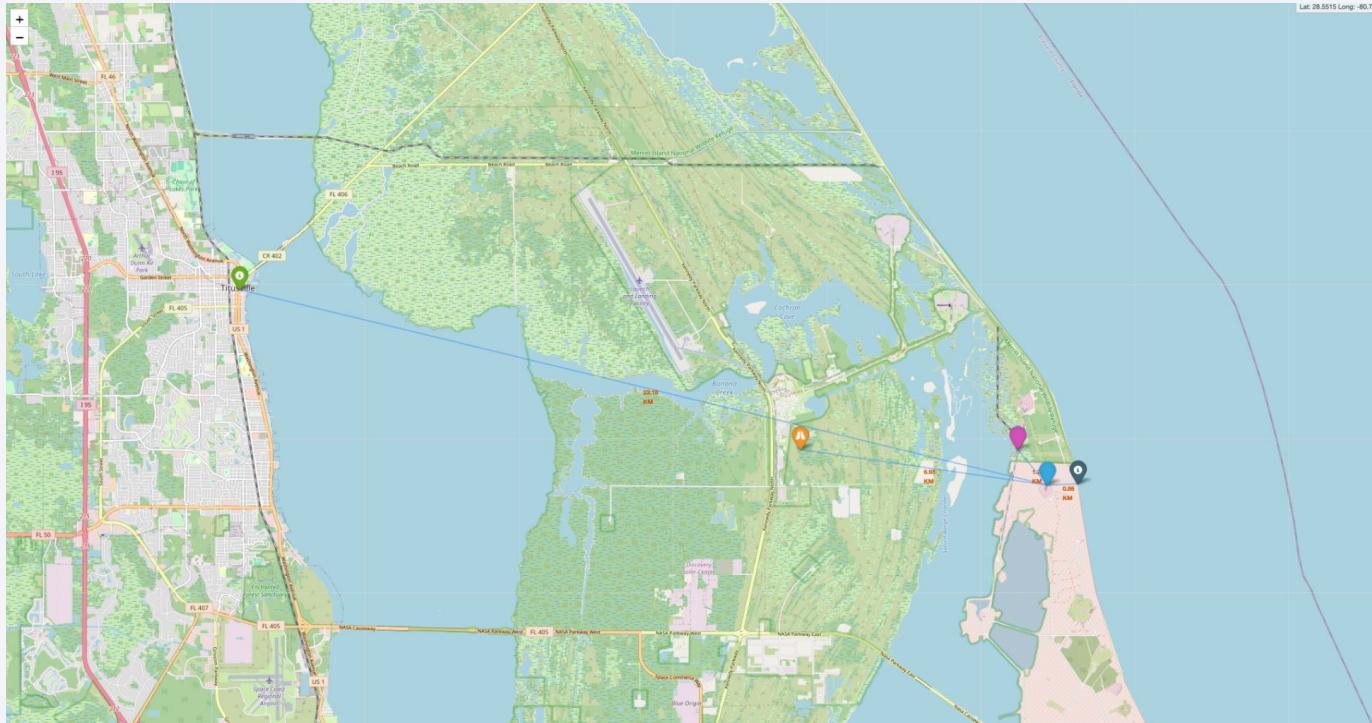


CCAFS SLC-40



1. These two launch sites, CCAFS LC-40 and CCAFS SLC-40, in Florida are in close proximity to each other as shown on the map
2. The individual launch sites are enclosed by their respective circles
3. The success (green) or failure (red) of a launch is annotated by individual mark ups for each launch site

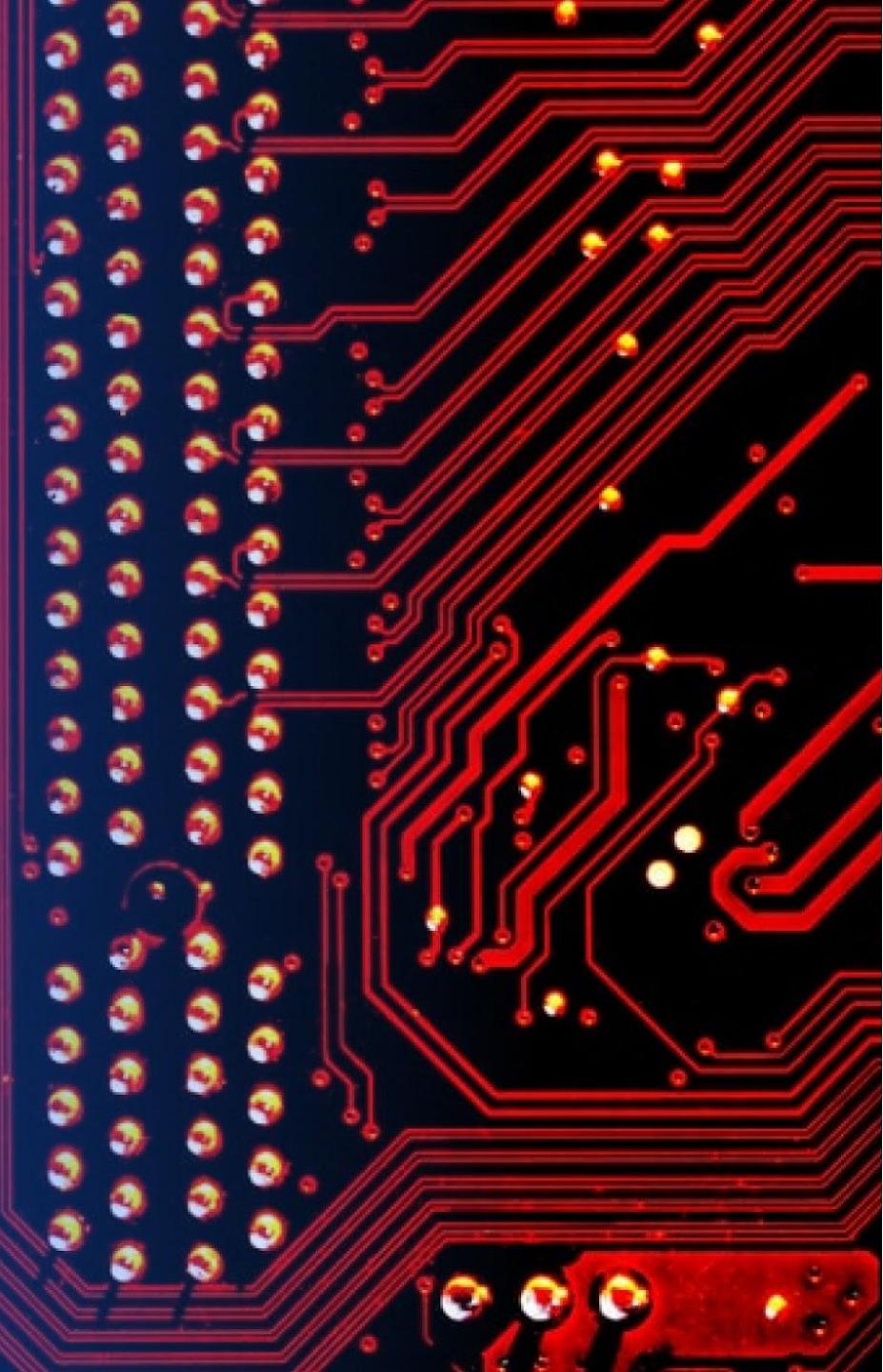
What is located near the launch site?



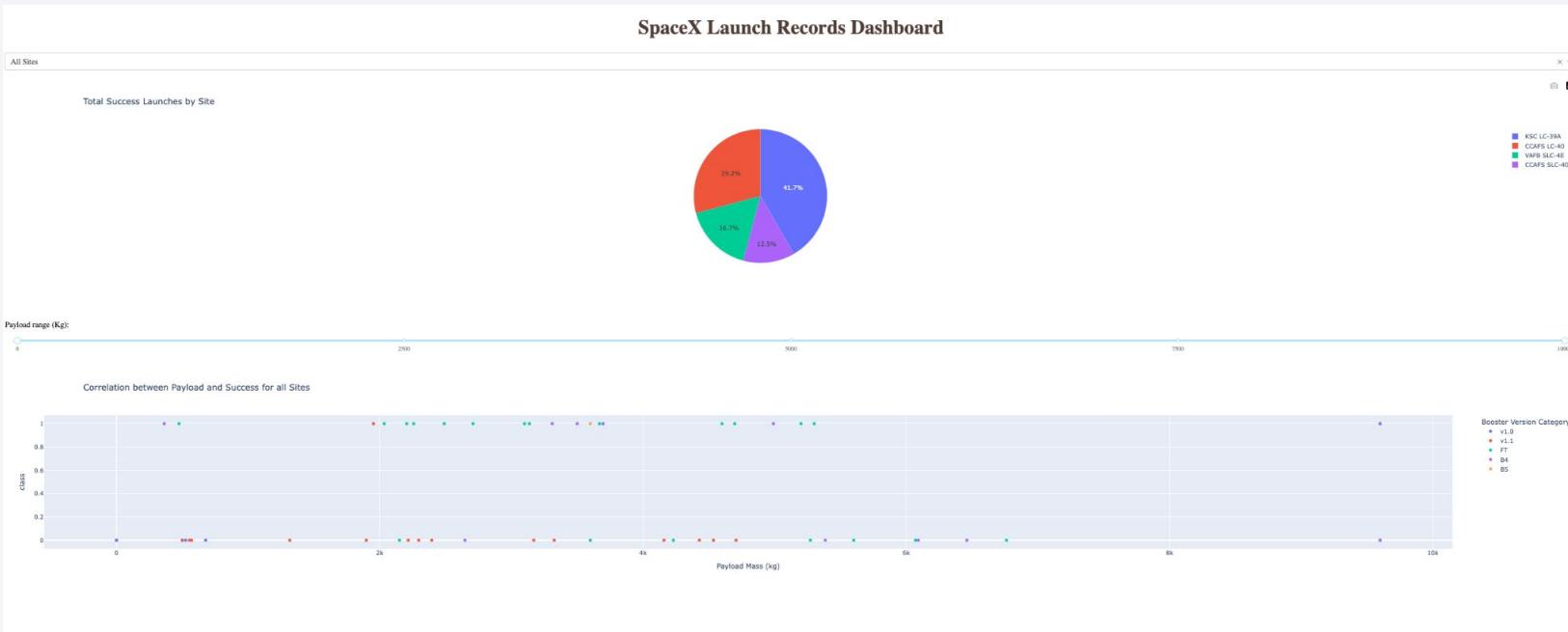
1. Proximity to railroads and highways
 - Strategic placement facilitates the transportation of heavy equipment and materials required for launches
 - Efficient logistics and operational Support
2. Proximity to coastlines
 - Launching over the ocean reduces risk to populated areas and allow for safer disposal of rocket stages and other debris
3. Distance from cities
 - Maintaining a safe distance is crucial for minimize risk to Potential human life and infrastructure
 - Mitigates potential hazards during launch operation

Section 4

Build a Dashboard with Plotly Dash



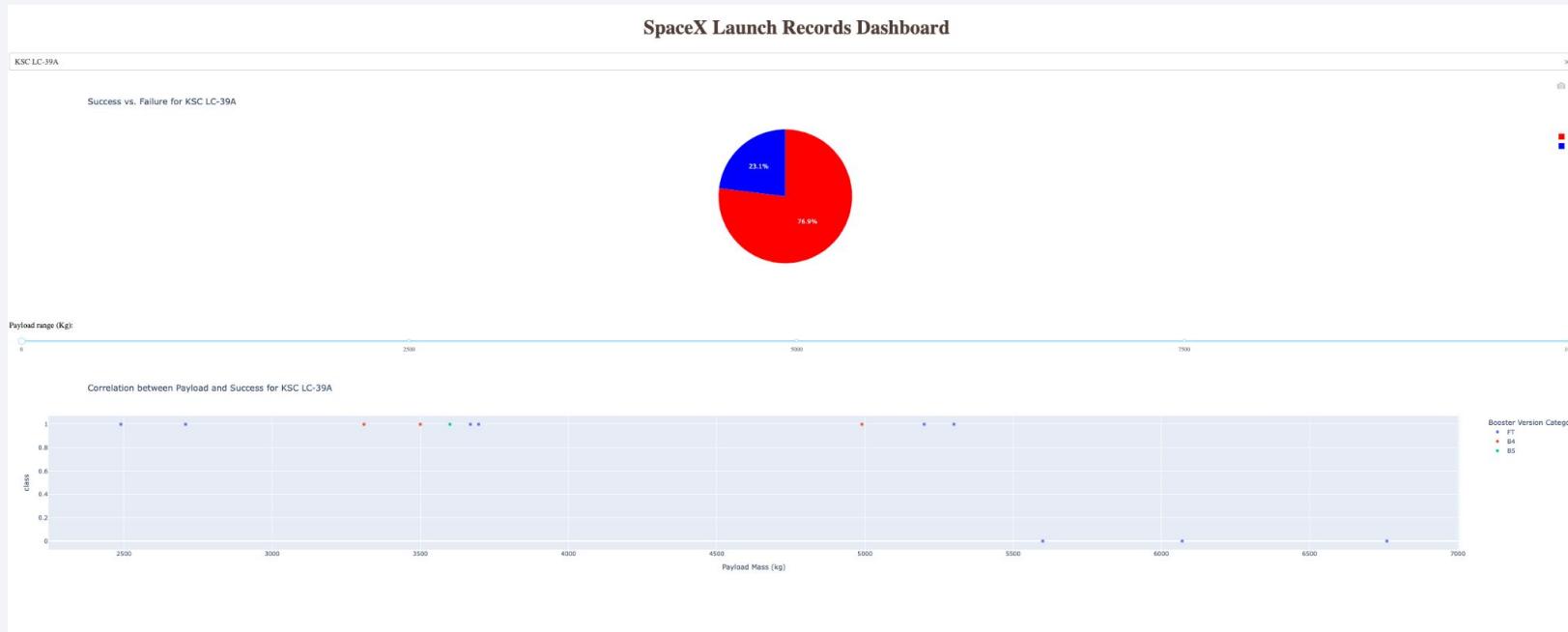
Successful Launches by Site



Dashboard Features

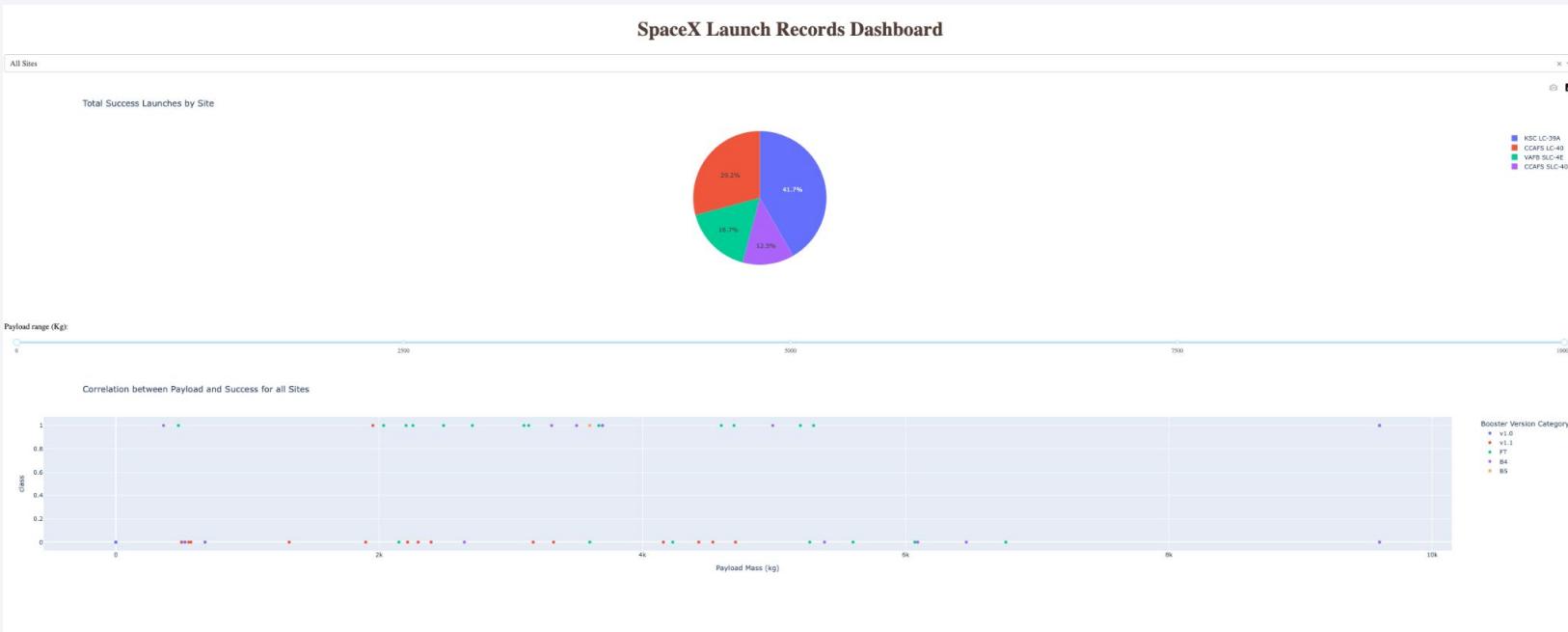
1. Range slider: Allows selection of payload mass that range from 0 kg up to 10,000 kg by increments of 2500 kg
2. Pie Chart callback: Updates based upon selected launch site
3. Scatter Chart callback: Updates scatter chart based upon selected launch site and payload range
 - a. Correlate launch success and payload mass to color-labeled dots representing booster version category

Highest Launch Success Ratio



1. KSC LC-39A has the highest launch success rate of 76.9% out of 13 attempts
2. Payloads to 5500 kg were successful whereas the three attempts at higher payloads (5500 - 7000 kg) failed
3. Makes up 41.7% of the total successful launches for SpaceX
4. Launches from KSC LC-39A contribute to 41.7% of the total successful launches for SpaceX

Successful Launch of Payloads

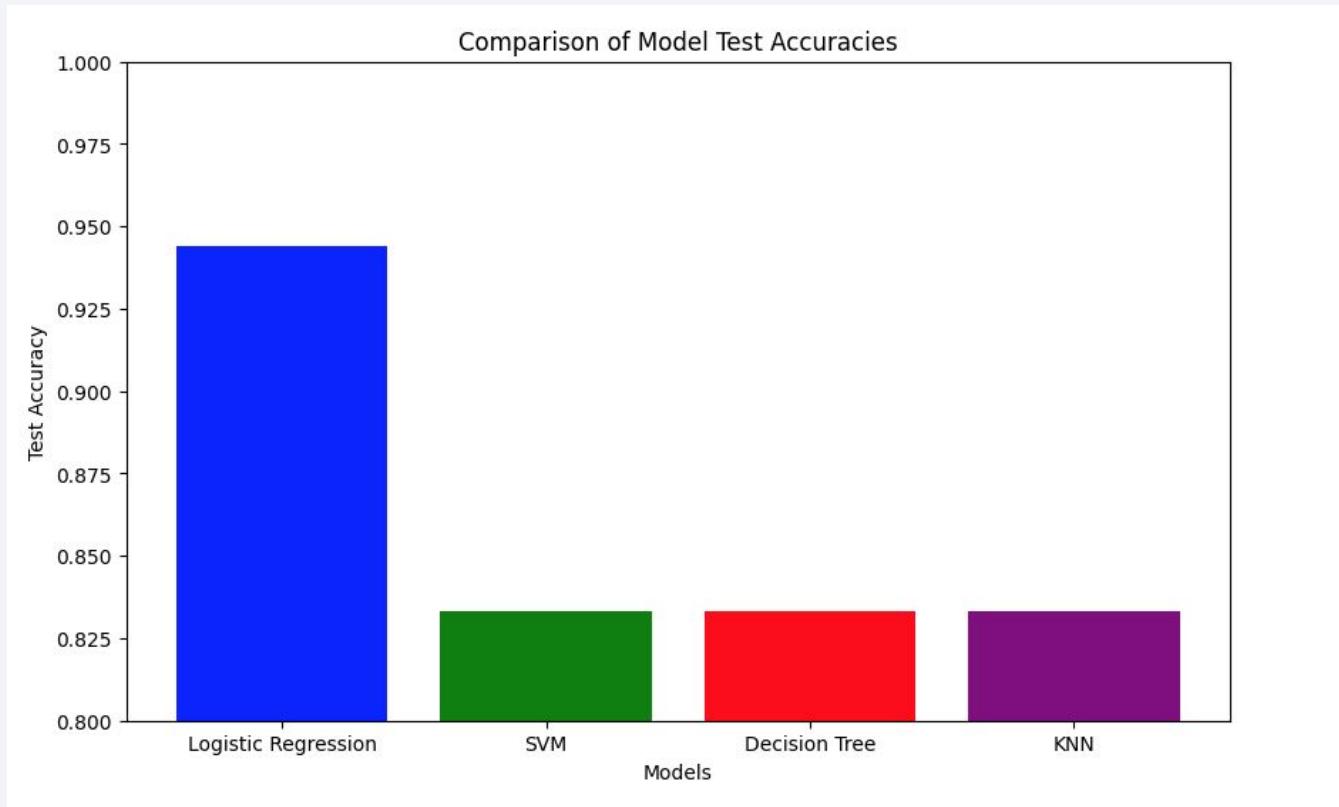


1. Highest rate of successful launches were with payloads in the 2000 - 4000 kg range
2. Least successful launches were with payloads greater than 6000 kg
3. Booster F9 FT was the most successful followed by F9 v1.1

Section 5

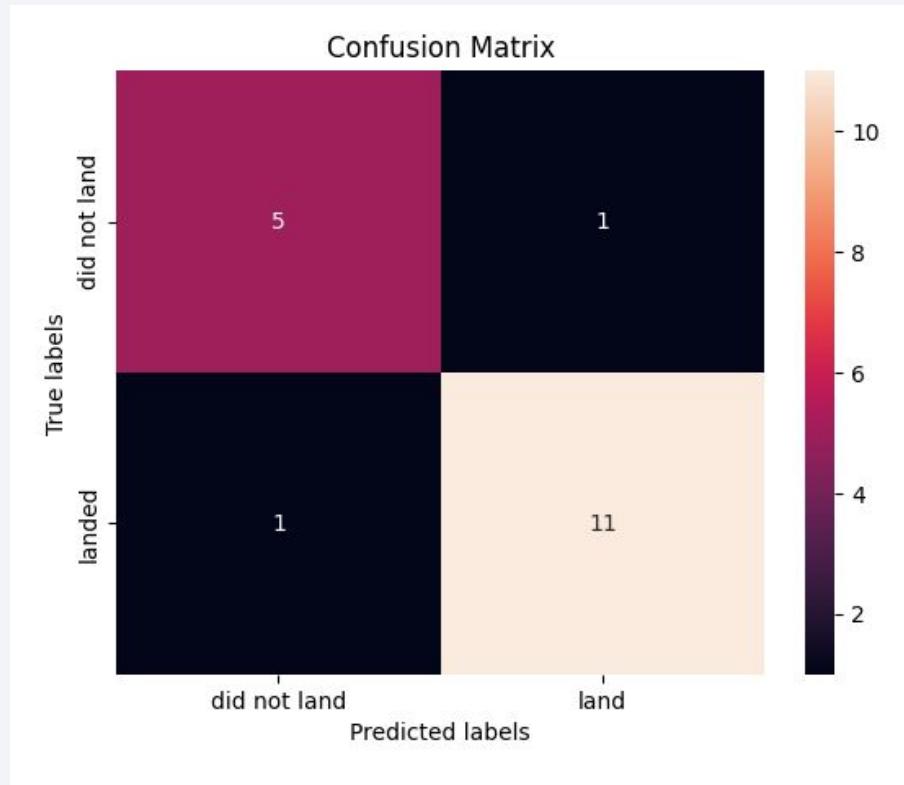
Predictive Analysis (Classification)

Classification Accuracy



1. Best accuracy for classification exhibited by Logistical Regression

Confusion Matrix



1. Logistical Regression demonstrated the highest accuracy (0.944) with fewer false positives and appears to be the most reliable based upon the confusion matrix
2. Logistical Regression is simpler and more interpretable

Conclusions

1. SpaceX was able to achieve a relative stable degree of success in their rocket launches through continually learning and innovating during the first 10 years of existence (2010-2020)
2. Although the frequency of recovery for the Falcon F9 first stage booster is positive, there is room for improvement in achieving certain orbits with the limited payload capacity
3. To be competitive, any newly formed entity needs to be conscientious that there exists a threshold of learning through execution that must be navigated in order to achieve success comparable to SpaceX.

Appendix

Appendix - Data Collection

```
# Define the URL
static_json_url = "https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBM-050321EN-SkillsNetwork/datasets/API_call_spacex_api.json"

# Perform the GET request to fetch the data
response = requests.get(static_json_url).json_normalize()

# Decode the response content as JSON
json_data = response.json()

# Convert the JSON response to a Pandas dataframe
data = pd.json_normalize(json_data)

# Take a subset of the dataframe
data = data[['rocket', 'payloads', 'launchpad', 'cores', 'flight_number', 'date_utc']]

# Remove rows with multiple cores and payloads
data = data[data['cores'].str.len() + data['payloads'].str.len() == 1]

# Extract the single value from lists
data['cores'] = data['cores'].map(lambda x: x[0])
data['payloads'] = data['payloads'].map(lambda x: x[0])

# Convert date_utc to datetime and extract the date
data['date'] = pd.to_datetime(data['date_utc']).dt.date

# Restrict the dates of the launches
import datetime
data = data[data['date'] >= datetime.date(2020, 11, 13)]

# Display the modified dataframe
print(data.head())

# rocket payloads \
0 Se90d0d95eda69955f7891eb Seb0e4b5b6c3b0006eebe1
1 Se90d0d95eda69955f7891eb Seb0e4b5b6c3b0006eebe2
3 Se90d0d95eda69955f7891eb Seb0e4b5b6c3b0006eebe1
4 Se90d0d95eda69955f7891eb Seb0e4b5b6c3b0006eebe16
5 Se90d0d95eda69973a0006eebe7

# launchpad \
0 Se9e450215090995d566f86
1 Se9e450215090995d566f86
3 Se9e450215090995d566f86
Se9e450215090995d566f86
5 Se9e4501f589094ba566f84

# cores \
0 {'core': 'Se9e289f35918033d3b2623', 'flight': 1, 'gridfins': False, 'legs': False, 'reused': False, 'landing_attempt': False, 'landing_success': None, 'landing_type': None, 'landpad': None}
1 {'core': 'Se9e289f35918416a3b2624', 'flight': 1, 'gridfins': False, 'legs': False, 'reused': False, 'landing_attempt': False, 'landing_success': None, 'landing_type': None, 'landpad': None}
3 {'core': 'Se9e289f35918416a3b2625', 'flight': 1, 'gridfins': False, 'legs': False, 'reused': False, 'landing_attempt': False, 'landing_success': None, 'landing_type': None, 'landpad': None}
4 {'core': 'Se9e289f35918418032627', 'flight': 1, 'gridfins': False, 'legs': False, 'reused': False, 'landing_attempt': False, 'landing_success': None, 'landing_type': None, 'landpad': None}
5 {'core': 'Se9e289f3591851203262628', 'flight': 1, 'gridfins': False, 'legs': False, 'reused': False, 'landing_attempt': False, 'landing_success': None, 'landing_type': None, 'landpad': None}

# flight_number date_utc date
0 1 2006-03-24T22:39:00.000Z 2006-03-24
1 2 2007-03-21T01:10:00.000Z 2007-03-21
3 4 2008-09-28T23:15:00.000Z 2008-09-28
4 5 2009-07-13T03:35:00.000Z 2009-07-13
5 6 2010-06-04T18:45:00.000Z 2010-06-06
```

Then, we need to create a Pandas data frame from the dictionary launch_dict.

```
[41]: # Create a data frame from launch_dict
import pandas as pd

# Sample dictionary containing the launch data
launch_dict = {
    'booster_version': 'BoosterVersion',
    'payload_mass': 'PayloadMass',
    'lat': 'Latitude',
    'long': 'Longitude',
    'latitude': 'Latitude',
    'longitude': 'Longitude',
    'landing_outcome': 'Outcome',
    'landing_type': 'LandingPad',
    'cores_reused': 'Reused',
    'cores_reused': 'Reused',
    'landing_pad': 'LandingPad',
    'core_block': 'CoreBlock',
    'times_reused': 'ReusedCount',
    'core_serial': 'Serial'
}

# Create a new Pandas dataframes from the dictionary
launch_df = pd.DataFrame(launch_dict)

# Display the first few rows of the new data frame
print(launch_df.head())

```

booster_version	payload_mass	lat	long	latitude	longitude	landing_outcome	landing_type	cores_reused	cores_reused	landing_pad	core_block	times_reused	core_serial
Falcon 1	28.0	None	None	None	None	None	None	False	False	None	None	1	False
Falcon 1	28.0	None	None	None	None	None	None	False	False	None	None	1	False
Falcon 1	28.0	None	None	None	None	None	None	False	False	None	None	1	False
Falcon 1	28.0	None	None	None	None	None	None	False	False	None	None	1	False

```
[41]: # Filter the data frame to only include Falcon 9 launches using .loc[]
data_falcon9 = launch_df.loc[launch_df['booster_version'] == 'Falcon 9']

# Print the first few rows of the filtered data frame
print(data_falcon9.head())

```

booster_version	payload_mass	lat	long	latitude	longitude	landing_outcome	landing_type	cores_reused	cores_reused	landing_pad	core_block	times_reused	core_serial
Falcon 9	28.0	None	None	None	None	None	None	False	False	None	None	1	False
Falcon 9	28.0	None	None	None	None	None	None	False	False	None	None	1	False
Falcon 9	28.0	None	None	None	None	None	None	False	False	None	None	1	False
Falcon 9	28.0	None	None	None	None	None	None	False	False	None	None	1	False
Falcon 9	28.0	None	None	None	None	None	None	False	False	None	None	1	False

Calculate below the mean for the PayloadMass using the .mean(). Then use the mean and the .replace()!

```
[41]: # Calculate the mean value of the PayloadMass column
# Replace the np.nan values with its mean value
import numpy as np

# Calculate the mean value of the PayloadMass column
mean_payload_mass = data_falcon9['payload_mass'].mean()

# Replace the np.nan values with the mean value using .loc[]
data_falcon9.loc[:, 'payload_mass'] = data_falcon9['payload_mass'].replace(np.nan, mean_payload_mass)

# Display the first few rows of the updated data frame
print(data_falcon9.head())

```

booster_version	payload_mass	lat	long	latitude	longitude	landing_outcome	landing_type	cores_reused	cores_reused	landing_pad	core_block	times_reused	core_serial
Falcon 9	6123.547647	None	None	None	None	None	None	False	False	None	None	1	False
Falcon 9	525.000000	None	None	None	None	None	None	False	False	None	None	1	False
Falcon 9	277.000000	ISS	CCSFS SLC 4E	-88.577396	-60.000000	ISS	CCSFS SLC 4E	-88.577396	-88.577396	ISS	CCSFS SLC 4E	-88.577396	False
Falcon 9	500.000000	PO	VAFB SLC 4E	-120.610829	317.000000	PO	VAFB SLC 4E	-120.610829	-120.610829	PO	VAFB SLC 4E	-120.610829	False
Falcon 9	3178.000000	GTO	CCSFS SLC 4E	-88.577396	None	GTO	CCSFS SLC 4E	-88.577396	-88.577396	GTO	CCSFS SLC 4E	-88.577396	False

latitude landing_outcome landing_type core_flight griffins_used \
4 28.561857 None None None 1 False
5 28.561857 None None None 1 False
6 28.561857 None None None 1 False
7 34.632093 False Ocean None 1 False
8 28.561857 None None None 1 False

core_reused legs_used landing_pad core_block times_reused core_serial \
4 False False None None 1 False
5 False False None None 1 False
6 False False None None 1 False
7 False False None None 1 False
8 False False None None 1 False

/home/jupyterlab/conda/envs/python/lib/python3.7/site-packages/pandas/core/indexing.py:1773: SettingWithCopyWarning: A value is trying to be set on a copy of a slice from a DataFrame. Try using .loc[row_indexer,col_indexer] = value instead.
See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#self-setters-single-columniloc[0,:], value, p1]
 booster_version \
0 None
payload_mass \
0 None
lat \
0 None
longitude \
0 None
latitude \
0 None
longitude \
0 None
landing_outcome \
0 None
landing_type \
26
core_reused \
0
griffins_used \
0
core_reused \
0
legs_used \
26
landing_pad \
26
core_block \
0
times_reused \
0
core_serial \
0
dtype: int64

You should see the number of missing values of the PayloadMass change to zero.

Appendix - Web Scraping

```

from bs4 import BeautifulSoup
import requests

# Define the URL for the Falcon 9 Wikipedia page
static_url = "https://en.wikipedia.org/wiki/Falcon_9"

# Perform the GET request to the URL
response = requests.get(static_url)

# Create a BeautifulSoup object from the response text content
soup = BeautifulSoup(response.text, 'html.parser')

# Find all table elements
html_tables = soup.find_all('table')

# Function to extract column name from header element
def extract_column_from_header(th):
    return th.text.strip()

# Initialize the list for column names
column_names = []

# Assuming the first table is the one we are interested in
first_launch_table = html_tables[0]

# 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
for th in first_launch_table.find_all('th'):
    name = extract_column_from_header(th)
    if name is not None and len(name) > 0:
        column_names.append(name)

# Display the extracted column names
print("Column Names: ", column_names)

# Create the dictionary from column names
launch_dict = dict.fromkeys(column_names)

# Remove an irrelevant column if it exists
if 'Date and time ()' in launch_dict:
    del launch_dict['Date and time ()']

# Initialize the launch_dict with each value to be an empty list
for key in launch_dict.keys():
    launch_dict[key] = []

# Add the new columns
new_columns = ['Function', 'Manufacturer', 'Country of origin', 'Cost per launch', 'Size', 'Height', 'Diameter', 'Mass d on', 'Derivatives', 'Launch history', 'Status', 'Launch sites', 'Total launches', 'Successes', 'Failure(s)', 'First flight', 'Last flight', 'Second flight', 'Diameter', 'Powered by', 'Maximum thrust', 'Specific impulse', 'Burn initial', 'launched', 'dict_keys():', 'Function', 'Manufacturer', 'Country of origin', 'Cost per launch', 'Size', 'rk', 'Launch history', 'Status', 'Launch sites', 'Total launches', 'Successes', 'Failure(s)', 'Partial failures', 'rsion Booster', 'Booster landing', 'Status', 'Time']]

# Display the initialized launch_dict keys to verify
print("Initialized launch_dict Keys: ", launch_dict.keys())

```

Appendix - Data Wrangling

TASK 1: Calculate the number of launches on each site

The data contains several Space X launch facilities: Cape Canaveral Space Launch Complex 40 VAFB SLC 4E, Vandenberg Air Force Base Space Launch Complex 4E (SLC-4E), Kennedy Space Center Launch Complex 39A KSC LC 39A .The location of each Launch is placed in the column `LaunchSite`.

Next, let's see the number of launches for each site.

Use the method `.value_counts()` on the column `LaunchSite` to determine the number of launches on each site:

```
# Print the column names of the DataFrame to verify the correct names
print(df.columns)

# Apply value_counts() on column LaunchSite
launch_site_counts = df['LaunchSite'].value_counts()

# Display the number of launches for each site
print(launch_site_counts)
```

Index(['FlightNumber', 'Date', 'BoosterVersion', 'PayloadMass', 'Orbit',
 'LaunchSite', 'Outcome', 'Flights', 'GridFins', 'Reused', 'Legs',
 'LandingPad', 'Block', 'ReusedCount', 'Serial', 'Longitude',
 'Latitude'],
 dtype='object')
LaunchSite
CCAFS SLC 40 55
KSC LC 39A 22
VAFB SLC 4E 13
Name: count, dtype: int64

Each launch aims to an dedicated orbit, and here are some common orbit types:

TASK 2: Calculate the number and occurrence of each orbit

Use the method `.value_counts()` to determine the number and occurrence of each orbit in the column `Orbit`:

```
# Print the column names of the DataFrame to verify the correct names
print(df.columns)

# Apply value_counts on Orbit column
orbit_counts = df['Orbit'].value_counts()

# Display the number and occurrence for each orbit
print(orbit_counts)
```

Index(['FlightNumber', 'Date', 'BoosterVersion', 'PayloadMass', 'Orbit',
 'LaunchSite', 'Outcome', 'Flights', 'GridFins', 'Reused', 'Legs',
 'LandingPad', 'Block', 'ReusedCount', 'Serial', 'Longitude',
 'Latitude'],
 dtype='object')
Orbit
GTO 27
ISS 21
VLEO 14
PO 9
LEO 7
SSO 5
MEO 3
HEO 1
ES-L1 1
SO 1
GEO 1
Name: count, dtype: int64

TASK 4: Create a landing outcome label from Outcome column

Using the `Outcome`, create a list where the element is zero if the corresponding row in `Outcome` is in the set `bad_outcome`; otherwise, it's one. Then assign it to the variable `landing_class`:

```
# landing_class = 0 if bad_outcome
# landing_class = 1 otherwise

# Define the set of bad outcomes
bad_outcome = {'False ASDS', 'False Ocean', 'False RTLS', 'None ASDS', 'None None'} # Adjust this set based on what constitutes a bad outcome in your data

# Create the landing_class list
landing_class = [0 if outcome in bad_outcome else 1 for outcome in df['Outcome']]

# Assign the landing_class list to the DataFrame
df['landing_class'] = landing_class

# Display the first few rows of the updated DataFrame
print(df.head())
```

FlightNumber	Date	BoosterVersion	PayloadMass	Orbit	LaunchSite	Outcome	Flights	GridFins	Reused	Legs	LandingPad	Block
0	2018-06-04	Falcon 9	6104.959412	LEO	CCAFS SLC 40	None None	1	False	False	False	NaN	1.0
1	2012-05-22	Falcon 9	525.000000	LEO	CCAFS SLC 40	None None	1	False	False	False	NaN	1.0
2	2013-03-01	Falcon 9	677.000000	ISS	CCAFS SLC 40	False Ocean	1	False	False	False	NaN	1.0
3	2013-09-29	Falcon 9	508.000000	PO	VAFB SLC 4E	None None	1	False	False	False	NaN	1.0
4	2013-12-03	Falcon 9	3178.000000	GTO	CCAFS SLC 40	None None	1	False	False	False	NaN	1.0

This variable will represent the classification variable that represents the outcome of each launch. If the value is zero, the first stage did not land successfully; one means the first stage landed Success!

TASK 3: Calculate the number and occurrence of mission outcome of the orbits

Use the method `.value_counts()` on the column `Outcome` to determine the number of `landing_outcomes`. Then assign it to a variable `landing_outcomes`:

```
[10]: # landing_outcomes = values on Outcome column
# Apply value_counts() on the 'landing_outcome' column to determine the number of landing_outcomes
landing_outcomes = df['Outcome'].value_counts()

# Display the number of landing outcomes
print(landing_outcomes)
```

Outcome
True ASDS 41
None None 19
True RTLS 14
False ASDS 6
True Ocean 5
False Ocean 2
None ASDS 2
False RTLS 1
Name: count, dtype: int64

True Ocean means the mission outcome was successfully landed to a specific region of the ocean while False Ocean means the mission outcome was unsuccessfully landed to a specific region of the ocean. True RTLS means the mission outcome was successfully landed to a ground pad. False RTLS means the mission outcome was unsuccessfully landed to a ground pad. True ASDS means the mission outcome was successfully landed to a drone ship. False ASDS means the mission outcome was unsuccessfully landed to a drone ship. None ASDS and None None: these represent a failure to land.

```
[11]: for i,outcome in enumerate(landing_outcomes.keys()):
    print(i,outcome)
```

0 True ASDS
1 None None
2 True RTLS
3 False ASDS
4 True Ocean
5 False Ocean
6 None ASDS
7 False RTLS

We create a set of outcomes where the second stage did not land successfully:

```
[12]: bad_outcomes=set([landing_outcomes.keys()[[1,3,5,6,7]]])
bad_outcomes
```

```
[12]: {'False ASDS', 'False Ocean', 'False RTLS', 'None ASDS', 'None None'}
```

Appendix - EDA Visual Analytics code

TASK 1: Visualize the relationship between Flight Number and Launch Site

Use the function `catplot` to plot `FlightNumber` vs `LaunchSite`, set the parameter `x` parameter to `FlightNumber`, set the `y` to `Launch Site` and set the parameter `hue` to `'class'`

```
# Plot a scatter point chart with x axis to be Flight Number and y axis to be the launch site, and hue to be the class value
import seaborn as sns
import matplotlib.pyplot as plt

# Plot a scatter point chart with x axis as Flight Number, y axis as Launch Site, and hue as class value
sns.catplot(x="FlightNumber", y="LaunchSite", hue="Class", data=df, aspect=3, kind="strip")

# Set the axis labels and title
plt.xlabel("Flight Number", fontsize=15)
plt.ylabel("Launch Site", fontsize=15)
plt.title("Flight Number vs. Launch Site", fontsize=20)

# Display the plot
plt.show()
```

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

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

```
: import seaborn as sns
import matplotlib.pyplot as plt

# Plot a scatter point chart with x axis as Payload Mass, y axis as Launch Site, and hue as class value
sns.catplot(x="PayloadMass", y="LaunchSite", hue="Class", data=df, aspect=3, kind="strip")

# Set the axis labels and title
plt.xlabel("Payload Mass (kg)", fontsize=15)
plt.ylabel("Launch Site", fontsize=15)
plt.title("Payload Mass vs. Launch Site", fontsize=20)

# Display the plot
plt.show()
```

TASK 3: Visualize the relationship between success rate of each orbit type

Next, we want to visually check if there are any relationship between success rate and orbit type.

Let's create a `bar chart` for the sucess rate of each orbit

```
# HINT use groupby method on Orbit column and get the mean of Class column

# Calculate the success rate for each orbit type
orbit_success_rate = df.groupby('Orbit')['Class'].mean().reset_index()

# Rename the columns for better understanding
orbit_success_rate.columns = ['Orbit', 'Success_Rate']

# Create a bar chart for the success rate of each orbit type
plt.figure(figsize=(12, 6))
plt.bar(orbit_success_rate['Orbit'], orbit_success_rate['Success_Rate'], color='skyblue')

# Set the axis labels and title
plt.xlabel("Orbit Type", fontsize=15)
plt.ylabel("Success Rate", fontsize=15)
plt.title("Success Rate by Orbit Type", fontsize=20)
plt.xticks(rotation=45)

# Display the plot
plt.show()
```

TASK 4: Visualize the relationship between FlightNumber and Orbit type

For each orbit, we want to see if there is any relationship between `FlightNumber` and `Orbit type`.

```
: # Plot a scatter point chart with x axis to be FlightNumber and y axis to be the Orbit, and hue to be the class value
import seaborn as sns
import matplotlib.pyplot as plt

# Plot a scatter point chart with x axis as Flight Number, y axis as Orbit, and hue as class value
sns.catplot(x="FlightNumber", y="Orbit", hue="Class", data=df, aspect=3, kind="strip")

# Set the axis labels and title
plt.xlabel("Flight Number", fontsize=15)
plt.ylabel("Orbit Type", fontsize=15)
plt.title("Flight Number vs. Orbit Type", fontsize=20)

# Display the plot
plt.show()
```

TASK 5: Visualize the relationship between Payload Mass and Orbit type

Similarly, we can plot the `Payload Mass` vs. `Orbit` scatter point charts to reveal the relationship between `Payload Mass` and `Orbit type`

```
: # Plot a scatter point chart with x axis to be Payload Mass and y axis to be the Orbit, and hue to be the class value
import seaborn as sns
import matplotlib.pyplot as plt

# Plot a scatter point chart with x axis as Payload Mass, y axis as Orbit, and hue as class value
sns.catplot(x="PayloadMass", y="Orbit", hue="Class", data=df, aspect=3, kind="strip")

# Set the axis labels and title
plt.xlabel("Payload Mass (kg)", fontsize=15)
plt.ylabel("Orbit Type", fontsize=15)
plt.title("Payload Mass vs. Orbit Type", fontsize=20)

# Display the plot
plt.show()
```

Appendix - EDA SQL

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

```
# Write the SQL query to get unique launch sites
query = 'SELECT DISTINCT "Launch_Site" FROM SPACEXTBL'

# Execute the query and load the results into a DataFrame
unique_launch_sites = pd.read_sql_query(query, conn)

# Display the results
print(unique_launch_sites)
```

	Launch_Site
0	CCAFS LC-40
1	VAFB SLC-4E
2	KSC LC-39A
3	CCAFS SLC-40

	Date	Time (UTC)	Booster_Version	Launch_Site
0	2010-06-04	18:45:00	F9 v1.0	B0003 CCAFS LC-40
1	2010-12-08	15:43:00	F9 v1.0	B0004 CCAFS LC-40
2	2012-05-22	7:44:00	F9 v1.0	B0005 CCAFS LC-40
3	2012-10-08	0:35:00	F9 v1.0	B0006 CCAFS LC-40
4	2013-03-01	15:10:00	F9 v1.0	B0007 CCAFS LC-40

	Payload	PAYOUTLOAD_MASS_KG_
0	Dragon Spacecraft Qualification Unit	0
1	Dragon demo flight C1, two CubeSats, barrel off...	0
2	Dragon demo flight C2	525
3	SpaceX CRS-1	500
4	SpaceX CRS-2	677

	Orbit	Customer	Mission_Outcome	Landing_Outcome
0	LEO	SpaceX	Success	Failure (parachute)
1	LEO (ISS)	NASA (COTS) NRO	Success	Failure (parachute)
2	LEO (ISS)	NASA (COTS)	Success	No attempt
3	LEO (ISS)	NASA (CRS)	Success	No attempt
4	LEO (ISS)	NASA (CRS)	Success	No attempt

Index(['Date', 'Time (UTC)', 'Booster_Version', 'Launch_Site', 'Payload', 'PAYLOAD_MASS_KG_', 'Orbit', 'Customer', 'Mission_Outcome', 'Landing_Outcome'],
 dtype='object')

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

```
# SQL query to get the average payload mass carried by booster version F9 v1.1
query = """
SELECT AVG("PAYLOAD_MASS_KG_") AS Average_Payload_Mass
FROM SPACEXTBL
WHERE "Booster_Version" = 'F9 v1.1';
"""

# Execute the query and load the results into a DataFrame
average_payload_mass = pd.read_sql_query(query, conn)

# Display the results
print(average_payload_mass)
```

	Average_Payload_Mass
0	2928.4

List the names of the booster_versions which have carried the maximum payload mass. Use a subquery

```
# SQL query to list booster versions with the maximum payload mass using a subquery
query = """
SELECT "Booster_Version"
FROM SPACEXTBL
WHERE "PAYLOAD_MASS_KG_" = (
    SELECT MAX("PAYLOAD_MASS_KG_")
    FROM SPACEXTBL
);
"""

# Execute the query and load the results into a DataFrame
boosters_max_payload = pd.read_sql_query(query, conn)

# Display the results
print(boosters_max_payload)
```

	Booster_Version
0	F9 B5 B1048.4
1	F9 B5 B1049.4
2	F9 B5 B1051.3
3	F9 B5 B1056.4
4	F9 B5 B1048.5
5	F9 B5 B1051.4
6	F9 B5 B1049.5
7	F9 B5 B1060.2
8	F9 B5 B1058.3
9	F9 B5 B1051.6
10	F9 B5 B1060.3
11	F9 B5 B1049.7

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.

```
# SQL query to rank the count of landing outcomes between the specified dates in descending order
query = """
SELECT "Landing_Outcome", COUNT(*) AS Outcome_Count
FROM SPACEXTBL
WHERE "Date" BETWEEN '2010-06-04' AND '2017-03-20'
GROUP BY "Landing_Outcome"
ORDER BY Outcome_Count DESC;
"""

# Execute the query and load the results into a DataFrame
landing_outcomes_ranked = pd.read_sql_query(query, conn)

# Display the results
print(landing_outcomes_ranked)
```

	Landing_Outcome	Outcome_Count
0	No attempt	10
1	Success (drone ship)	5
2	Failure (drone ship)	5
3	Success (ground pad)	3
4	Controlled (ocean)	3
5	Uncontrolled (ocean)	2
6	Failure (parachute)	2
7	Precluded (drone ship)	1

Appendix - Visualizations

```

# Create a marker with distance to a closest city, railway, highway, etc.
# Draw a line between the marker to the launch site

import folium
from folium.plugins import MarkerCluster, MousePosition
from folium.features import DivIcon
import numpy as np

def haversine(lat1, lon1, lat2, lon2):
    # Convert latitude and longitude from degrees to radians
    lat1, lon1, lat2, lon2 = map(np.radians, [lat1, lon1, lat2, lon2])

    # Haversine formula
    dlat = lat2 - lat1
    dlon = lon2 - lon1
    a = np.sin(dlat/2)**2 + np.cos(lat1) * np.cos(lat2) * np.sin(dlon/2)**2
    c = 2 * np.arcsin(np.sqrt(a))
    r = 6371 # Radius of Earth in kilometers. Use 3956 for miles.
    return c * r

# Coordinates of the launch site (CCAFS SLC-40)
launch_site_lat = 28.563197
launch_site_lon = -80.576820

# Coordinates of the closest coastline point
coastline_lat = 28.56353
coastline_lon = -80.56798

# Coordinates of Titusville, FL
titusville_lat = 28.6122
titusville_lon = -80.8076

# Coordinates of the FEC Railway near Cape Canaveral, FL
railway_lat = 28.5721
railway_lon = -80.5953

# Coordinates of the nearest point on Max Brewer Causeway (Route 405)
highway_lat = 28.5721
highway_lon = -80.6473

# Calculate distances
distance_coastline = haversine(launch_site_lat, launch_site_lon, coastline_lat, coastline_lon)
distance_titusville = haversine(launch_site_lat, launch_site_lon, titusville_lat, titusville_lon)
distance_railway = haversine(launch_site_lat, launch_site_lon, railway_lat, railway_lon)
distance_highway = haversine(launch_site_lat, launch_site_lon, highway_lat, highway_lon)

print("Distance between the launch site and the closest coastline point: {distance_coastline:.2f} km")
print("Distance between the launch site and Titusville, FL: {distance_titusville:.2f} km")
print("Distance between the launch site and the FEC Railway: {distance_railway:.2f} km")
print("Distance between the launch site and the nearest highway: {distance_highway:.2f} km")

# Start location is NASA Johnson Space Center
nasa_coordinate = [29.55964888503615, -95.0830971930759]
site_map = folium.Map(location=nasa_coordinate, zoom_start=5)

# Initialize MarkerCluster
marker_cluster = MarkerCluster().add_to(site_map)

# Add a marker for the launch site (CCAFS SLC-40)
launch_site_marker = folium.Marker(
    location=[launch_site_lat, launch_site_lon],
    popup="CCAFS SLC-40 Launch Site",
    icon=folium.Icon(color='blue', icon='rocket')
)
site_map.add_child(launch_site_marker)

# Add a marker for the closest coastline point
coastline_marker = folium.Marker(
    location=[coastline_lat, coastline_lon],
    popup="Closest Coastline",
    icon=folium.Icon(color='cadetblue', icon='info-sign')
)
site_map.add_child(coastline_marker)

# Add a marker for Titusville, FL
titusville_marker = folium.Marker(
    location=[titusville_lat, titusville_lon],
    popup="Titusville, FL",
    icon=folium.Icon(color='green', icon='info-sign')
)
site_map.add_child(titusville_marker)

# Add a marker for the FEC Railway
railway_marker = folium.Marker(
    location=[railway_lat, railway_lon],
    popup="FEC Railway",
    icon=folium.Icon(color='purple', icon='train')
)
site_map.add_child(railway_marker)

# Add a marker for the nearest highway (Max Brewer Causeway)
highway_marker = folium.Marker(
    location=[highway_lat, highway_lon],
    popup="Nearest Highway (Max Brewer Causeway)",
    icon=folium.Icon(color='orange', icon='road')
)
site_map.add_child(highway_marker)

# Add markers displaying the distances
distance_markers = [
    ((launch_site_lat + coastline_lat) / 2, (launch_site_lon + coastline_lon) / 2, distance_coastline, 'coastline', '#d35400'),
    ((launch_site_lat + titusville_lat) / 2, (launch_site_lon + titusville_lon) / 2, distance_titusville, 'Titusville', '#d35400'),
    ((launch_site_lat + railway_lat) / 2, (launch_site_lon + railway_lon) / 2, distance_railway, 'railway', '#d35400'),
    ((launch_site_lat + highway_lat) / 2, (launch_site_lon + highway_lon) / 2, distance_highway, 'highway', '#d35400')
]

for lat, lon, distance, label, color in distance_markers:
    distance_marker = folium.Marker(
        location=[lat, lon],
        icon=DivIcon(
            icon_size=(20, 20),
            icon_anchor=(0, 0),
            html=f"<div style='font-size: 12; color:{color};'><b>{distance:.2f}</b> KM</div>"
        )
    )
    site_map.add_child(distance_marker)

# Create PolyLines between the launch site and each location
coordinates_list = [
    ((launch_site_lat, launch_site_lon), (coastline_lat, coastline_lon)),
    ((launch_site_lat, launch_site_lon), (titusville_lat, titusville_lon)),
    ((launch_site_lat, launch_site_lon), (railway_lat, railway_lon)),
    ((launch_site_lat, launch_site_lon), (highway_lat, highway_lon))
]

for coordinates in coordinates_list:
    lines = folium.PolyLine(locations=coordinates, weight=1)
    site_map.add_child(lines)

```

```

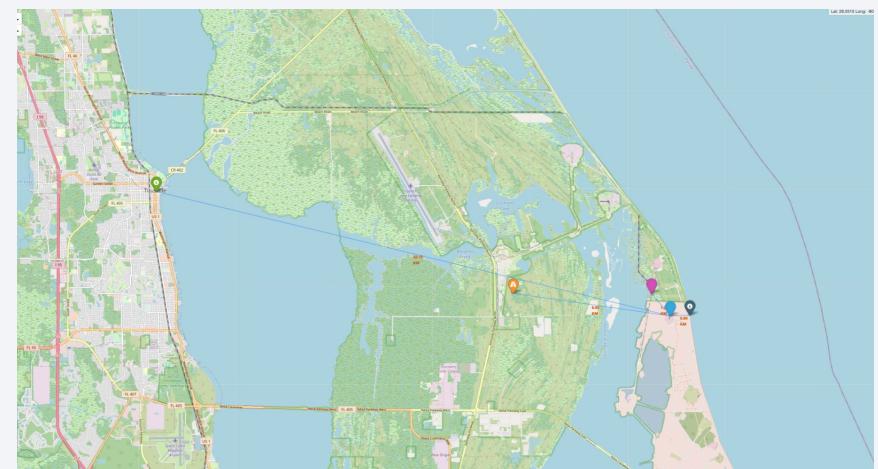
# Add Mouse Position to get the coordinate (Lat, Long) for a mouse over on the map
formatter = "function(num) {return L.Util.formatNum(num, 5);}"
mouse_position = MousePosition(
    position='topright',
    separator=' Long: ',
    empty_string='NaN',
    lng_first=False,
    num_digits=5,
    prefix='Lat:',
    lat_formatter=formatter,
    lng_formatter=formatter,
)

site_map.add_child(mouse_position)

# Display the map
site_map

Distance between the launch site and the closest coastline point: 0.86 km
Distance between the launch site and Titusville, FL: 23.18 km
Distance between the launch site and the FEC Railway: 1.29 km
Distance between the launch site and the nearest highway: 6.95 km

```



Appendix - Dashboard

```

import dash
from dash import dcc, html, Input, Output
import pandas as pd
import plotly.express as px

# Read the SpaceX launch data into pandas dataframe
spacex_df = pd.read_csv("spacex_launch_dash.csv")
max_payload = spacex_df['Payload Mass (kg)'].max()
min_payload = spacex_df['Payload Mass (kg)'].min()

# Create a dash application
app = dash.Dash(__name__)

# Create an app layout
app.layout = html.Div(children=[

    html.H1('SpaceX Launch Records Dashboard',
        style={'textAlign': 'center', 'color': '#503D36',
               'font-size': 40}),

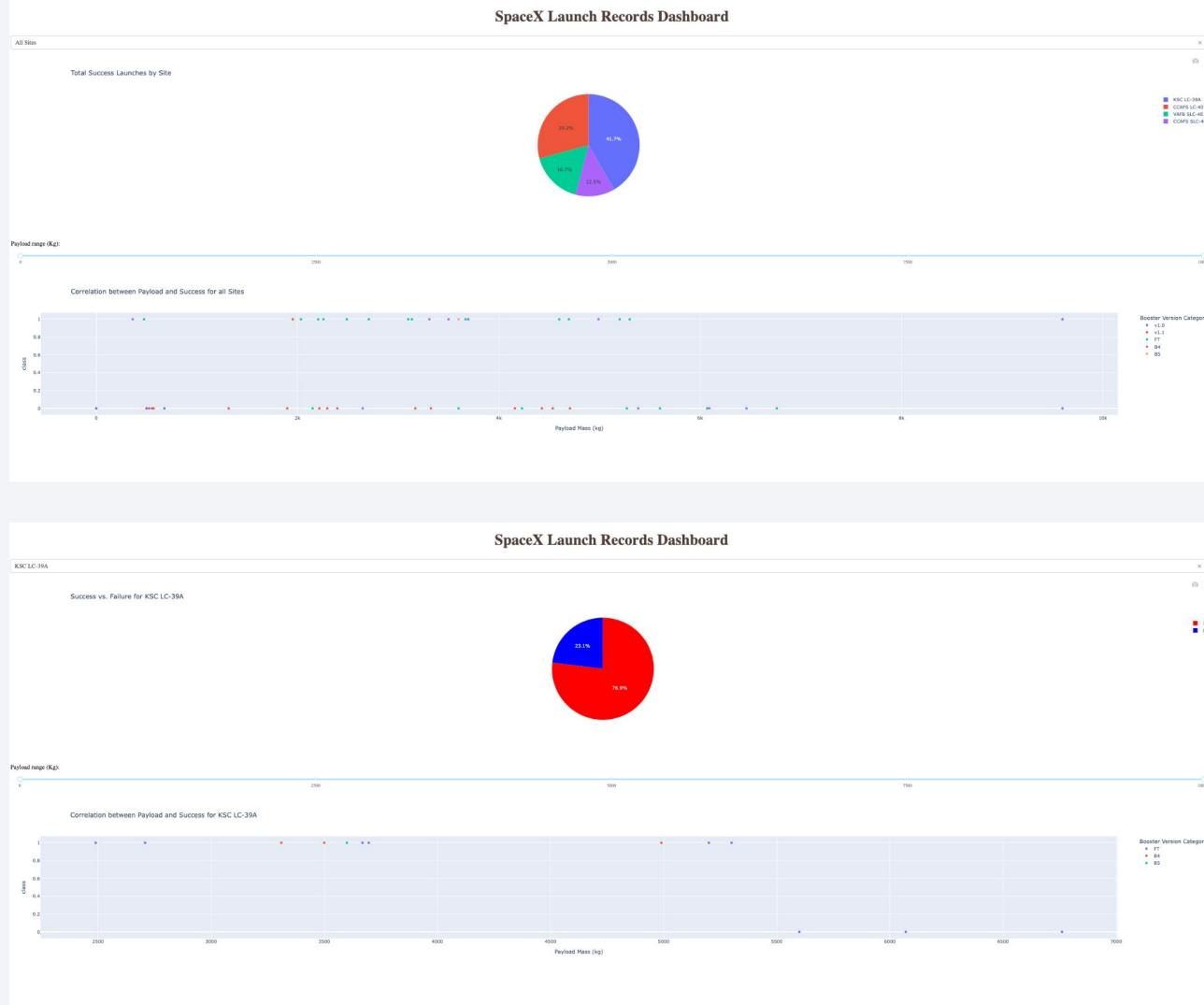
    dcc.Dropdown(id='site-dropdown',
        options=[
            {'label': 'All Sites', 'value': 'ALL'},
            {'label': 'CCAFS LC-40', 'value': 'CCAFS LC-40'},
            {'label': 'CCAFS SLC-40', 'value': 'CCAFS SLC-40'},
            {'label': 'KSC LC-39A', 'value': 'KSC LC-39A'},
            {'label': 'VAFB SLC-4E', 'value': 'VAFB SLC-4E'}
        ],
        value='ALL',
        placeholder="Select a Launch Site here",
        searchable=True),
    html.Br(),
    html.Div(dcc.Graph(id='success-pie-chart')),
    html.Br(),
    html.P("Payload range (Kg):"),
    dcc.RangeSlider(id='payload-slider',
        min=0, max=10000, step=1000,
        marks={0: '0', 2500: '2500', 5000: '5000', 7500: '7500', 10000: '10000'},
        value=[min_payload, max_payload]),
    html.Br(),
    html.Div(dcc.Graph(id='success-payload-scatter-chart')),
])

@app.callback(
    Output(component_id='success-pie-chart', component_property='figure'),
    Input(component_id='site-dropdown', component_property='value')
)
def get_pie_chart(entered_site):
    color_map = {0: 'blue', 1: 'red'} # Specify the color map for class
    if entered_site == 'ALL':
        success_counts = spacex_df[spacex_df['class'] == 1].groupby(['Launch Site']).size().reset_index(name='count')
        fig = px.pie(success_counts,
                     values='count',
                     names='Launch Site',
                     title='Total Success Launches by Site')
    else:
        filtered_df = spacex_df[spacex_df['Launch Site'] == entered_site]
        success_counts = filtered_df.groupby('class').size().reset_index(name='count')
        fig = px.pie(success_counts,
                     values='count',
                     names='class',
                     title=f"Success vs. Failure for {entered_site}",
                     color='class',
                     color_discrete_map=color_map)
    return fig

@app.callback(
    Output(component_id='success-payload-scatter-chart', component_property='figure'),
    [Input(component_id='site-dropdown', component_property='value'),
     Input(component_id='payload-slider', component_property='value')]
)
def update_scatter(entered_site, payload_range):
    low, high = payload_range
    filtered_df = spacex_df[(spacex_df['Payload Mass (kg)'] >= low) & (spacex_df['Payload Mass (kg)'] <= high)]
    if entered_site == 'ALL':
        fig = px.scatter(filtered_df,
                         x='Payload Mass (kg)',
                         y='class',
                         color='Booster Version Category',
                         title='Correlation between Payload and Success for all Sites')
    else:
        filtered_df = filtered_df[filtered_df['Launch Site'] == entered_site]
        fig = px.scatter(filtered_df,
                         x='Payload Mass (kg)',
                         y='class',
                         color='Booster Version Category',
                         title=f'Correlation between Payload and Success for {entered_site}')
    return fig

if __name__ == '__main__':
    app.run_server(debug=True)

```



Appendix - Machine Learning Prediction

```
[ ] # Modified Grid Search
import numpy as np
from sklearn.model_selection import RandomizedSearchCV
from sklearn.svm import SVC

# Define a simplified parameter grid
parameters = {'kernel': ('linear', 'rbf'),
              'C': [0.1, 1, 10],
              'gamma': [0.001, 0.01]}

# Create the SVM object
svm = SVC()

# Create the RandomizedSearchCV object with cv = 10 and parallel computation
svm_cv = RandomizedSearchCV(svm, parameters, n_iter=10, cv=10, n_jobs=-1)

# Fit the RandomizedSearchCV object to find the best parameters
svm_cv.fit(X_train, Y_train)

# Print the best parameters and the accuracy
print("Tuned hyperparameters: (best parameters) ", svm_cv.best_params_)
print("Accuracy: ", svm_cv.best_score_)

# Calculate the accuracy on the test data using the best estimator
test_accuracy = svm_cv.best_estimator_.score(X_test, Y_test)
print("Test Accuracy: ", test_accuracy)
```

```
[ ] # Initial Data Tree
import numpy as np
from sklearn.model_selection import GridSearchCV
from sklearn.tree import DecisionTreeClassifier
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
import pandas as pd

# Load your data from local files
data = pd.read_csv("./content/dataset_part_2.csv")
X = pd.read_csv("./content/dataset_part_3.csv")

# Convert 'Class' column to a NumPy array
Y = data['Class'].to_numpy()

# Standardize the data
scaler = StandardScaler()
X = scaler.fit_transform(X)

# Split the data into training and test sets
X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.2, random_state=2)

# Define the parameter grid
parameters = {'criterion': ['gini', 'entropy'],
              'splitter': ['best', 'random'],
              'max_depth': [2*n for n in range(1,10)],
              'max_features': ['auto', 'sqrt'],
              'min_samples_leaf': [1, 2, 4],
              'min_samples_split': [2, 5, 10]}

# Create the Decision Tree Classifier object
tree = DecisionTreeClassifier()
```

```
[ ] from sklearn.model_selection import GridSearchCV
from sklearn.linear_model import LogisticRegression

# Define parameter grid with different solvers and increased max_iter
param_grid = {'C': [0.1, 1, 10, 100], 'solver': ['liblinear']}
logreg = LogisticRegression(max_iter=2000) # Increase max_iter to 2000

# Initialize GridSearchCV
logreg_cv = GridSearchCV(logreg, param_grid, cv=5)

# Assuming X_train and Y_train are already defined and preprocessed
logreg_cv.fit(X_train, Y_train)

# Print the best parameters and the accuracy
print("Tuned hyperparameters: (best parameters) ", logreg_cv.best_params_)
print("Accuracy: ", logreg_cv.best_score_)
```

```
[ ] # Initial KNN
import numpy as np
from sklearn.model_selection import GridSearchCV, train_test_split
from sklearn.neighbors import KNeighborsClassifier
from sklearn.preprocessing import StandardScaler
import pandas as pd

# Load your data from local files
data = pd.read_csv("./content/dataset_part_2.csv")
X = pd.read_csv("./content/dataset_part_3.csv")

# Convert 'Class' column to a NumPy array
Y = data['Class'].to_numpy()

# Standardize the data
scaler = StandardScaler()
X = scaler.fit_transform(X)

# Split the data into training and test sets
X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.2, random_state=2)

# Define the parameter grid
parameters = {'n_neighbors': [1, 2, 3, 4, 5, 6, 7, 8, 9, 10],
              'algorithm': ['auto', 'ball_tree', 'kd_tree', 'brute'],
              'p': [1, 2]}

# Create the KNN Classifier object
KNN = KNeighborsClassifier()
```

Thank you!

