



IBM Developer  
SKILLS NETWORK

# Winning Space Race with Data Science

Liao Chang  
11/08/2023



# Outline

---

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

# Executive Summary

---

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

# Introduction

---

- SpaceX, a leader in the space industry, strives to make space travel affordable for everyone. Its accomplishments include sending spacecraft to the international space station, launching a satellite constellation that provides internet access and sending manned missions to space. SpaceX can do this because the rocket launches are relatively inexpensive (\$62 million per launch) due to its novel reuse of the first stage of its Falcon 9 rocket. Other providers, which are not able to reuse the first stage, cost upwards of \$165 million each. By determining if the first stage will land, we can determine the price of the launch. To do this, we can use public data and machine learning models to predict whether SpaceX – or a competing company – can reuse the first stage.



Section 1

# Methodology



# Methodology

---

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

# Data Collection

---

- Data are collected via
  - Web scraping SpaceX Wikipedia page
  - SpaceX REST API

# Data Collection – SpaceX API

---

- [Github](#)

- **Request data** from SpaceX API (rocket launch data)
- **Decode response** using `.json()` and convert to a dataframe using `.json_normalize()`
- **Request information** about the launches from SpaceX API using custom functions
- **Create dictionary** from the data
- **Create dataframe** from the dictionary
- **Filter dataframe** to contain only Falcon 9 launches
- **Replace missing values** of Payload Mass with calculated `.mean()`
- **Export data** to csv file



# Data Collection - Scraping

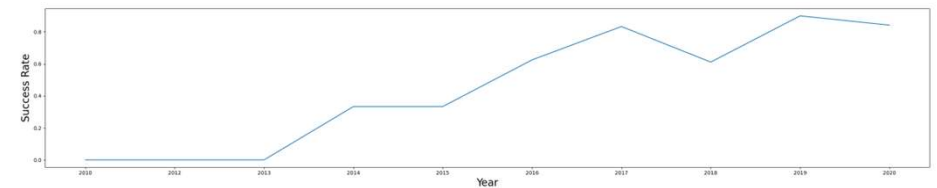
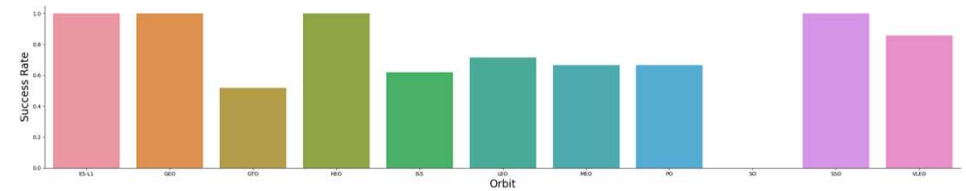
---

- [Github Link](#)

- **Request data** (Falcon 9 launch data) from Wikipedia
- **Create BeautifulSoup object** from HTML response
- **Extract column names** from HTML table header
- **Collect data** from parsing HTML tables
- **Create dictionary** from the data
- **Create dataframe** from the dictionary
- **Export data** to csv file

# EDA with Data Visualization

- [Github link](#)
- **Create charts** to analyze relationships and show comparisons



# EDA with SQL

- [Github link](#)
- Query the data to understand more about the data

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

## Task 10

Rank the count of landing outcomes (such as Failure (drone ship) or Success (ground pad)) between the date 2010-06-04 and 2017-03-20, in descending order.

```
0]: %sql select landing_outcome, count(*) as cnt from spacetable group by landing_outcome having date between '2010-06-04' and '2017-03-20'
```

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

```
0]:
```

Landing_Outcome	cnt
No attempt	21
Success (drone ship)	14
Success (ground pad)	9
Failure (drone ship)	5
Controlled (ocean)	5
Uncontrolled (ocean)	2
Precluded (drone ship)	1

# Build an Interactive Map with Folium

---

- [Github Link](#)
- **Create maps** to visualize launch sites, view launch outcomes and see distance to proximities
-

# Build a Dashboard with Plotly Dash

---

- **Create dashboard**
- Pie chart showing successful launches
- Scatter chart showing Payload Mass vs. Success Rate by Booster Version
- [Link](#)

# Predictive Analysis (Classification)

---

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

# Results

---

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



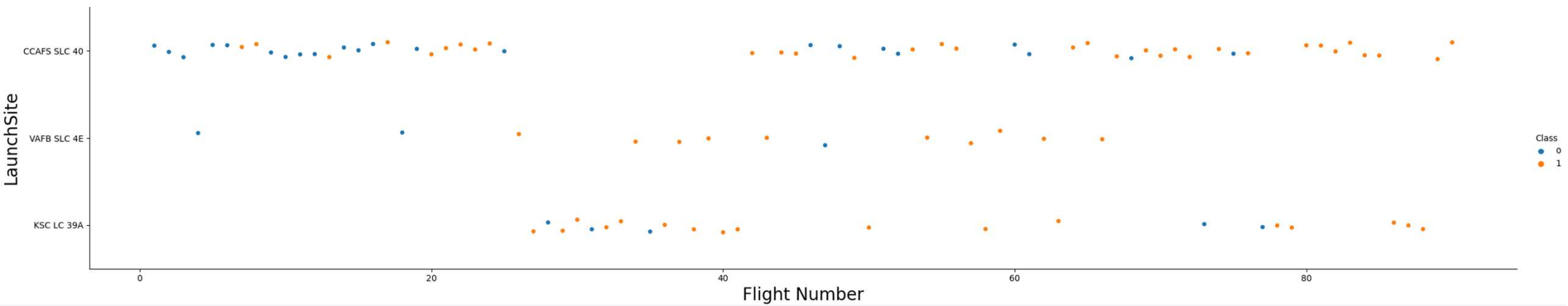


Section 2

# Insights drawn from EDA

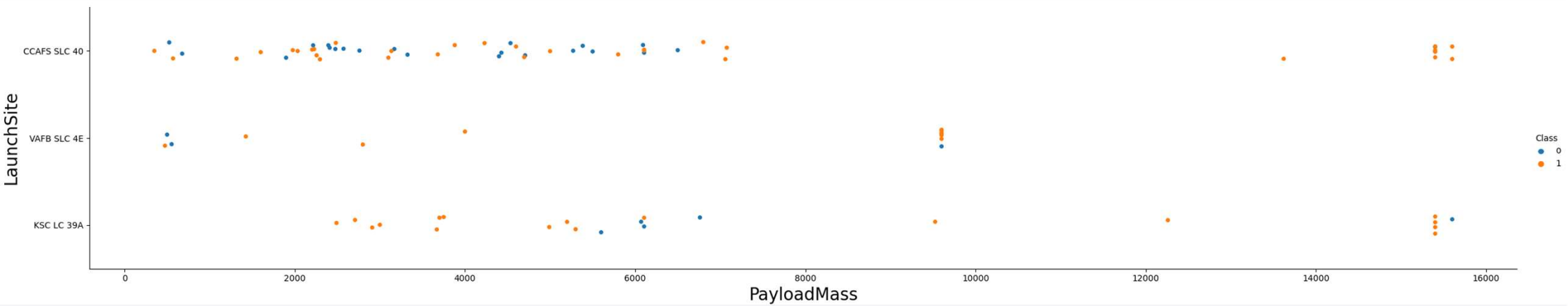
# Flight Number vs. Launch Site

---



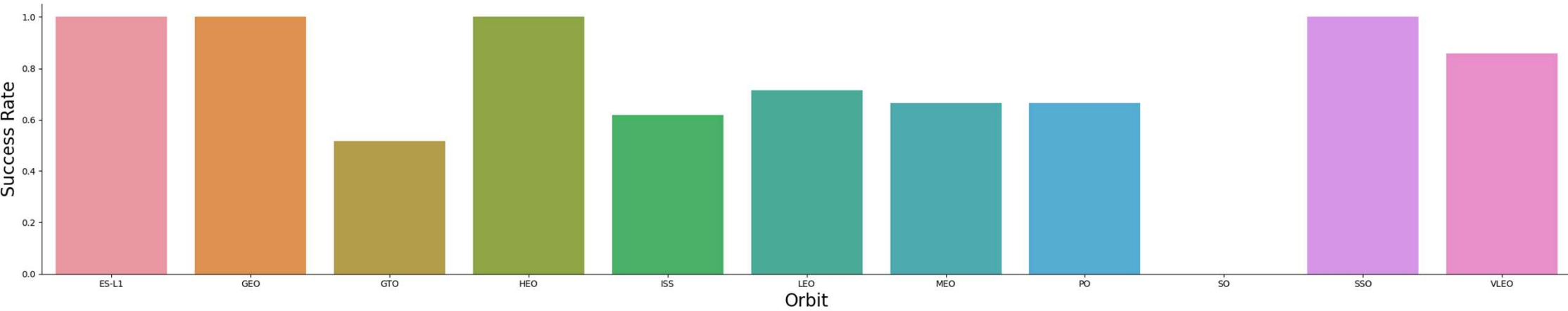
# Payload vs. Launch Site

---



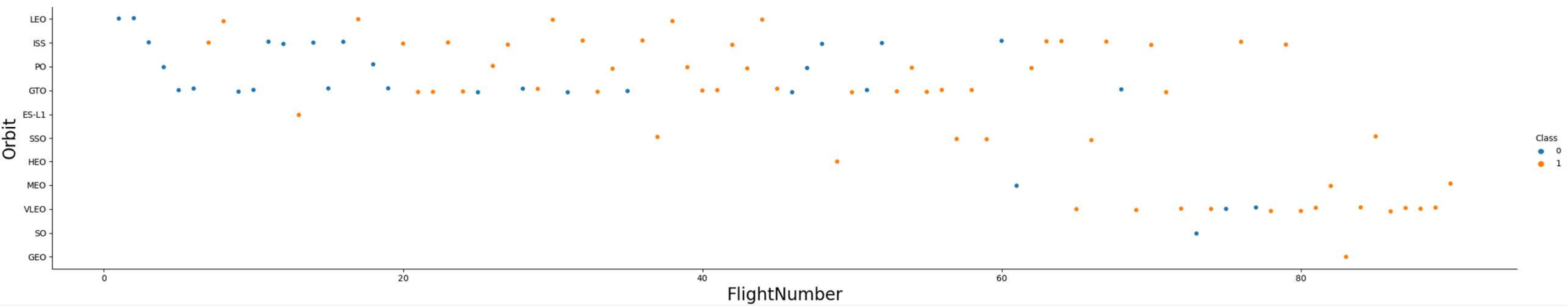
# Success Rate vs. Orbit Type

---



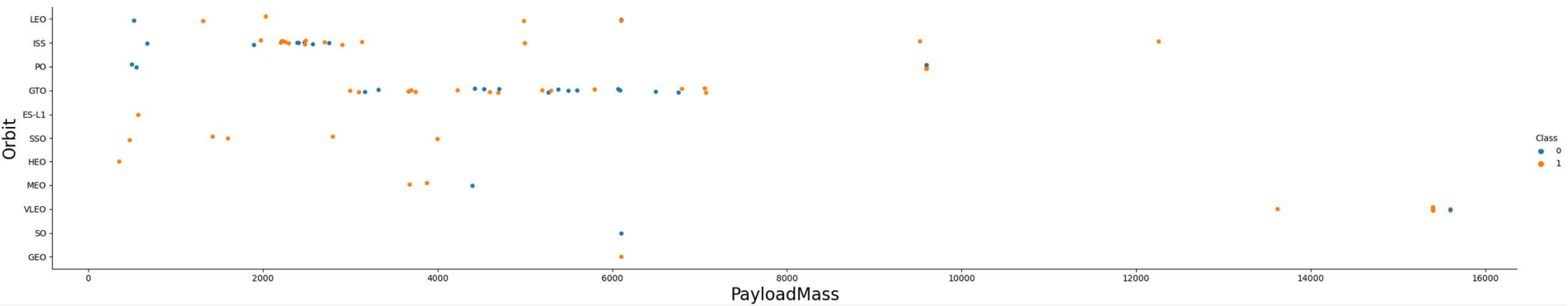
# Flight Number vs. Orbit Type

---



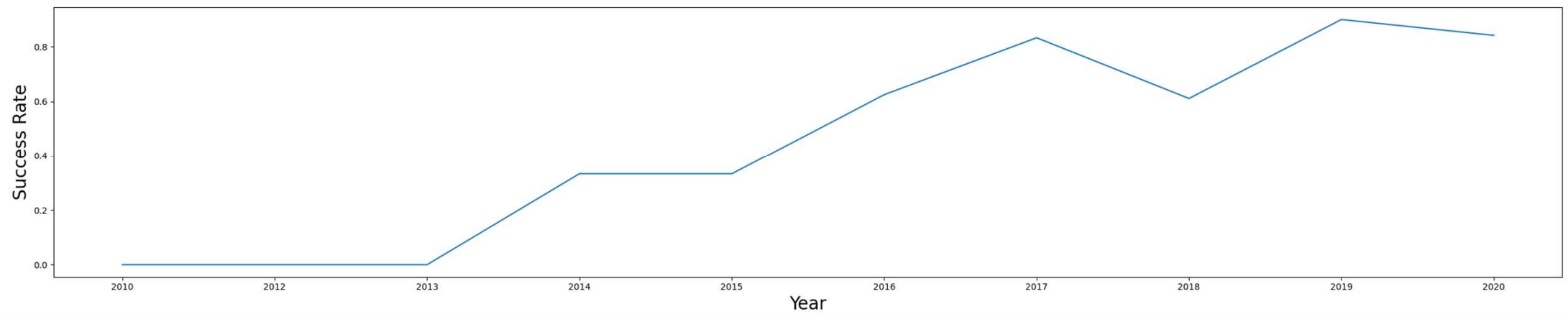
# Payload vs. Orbit Type

---



# Launch Success Yearly Trend

---





# All Launch Site Names

---

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

# Launch Site Names Begin with 'CCA'

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

\* sqlite:///my\_data1.db  
Done.

Date	Time (UTC)	Booster_Version	Launch_Site	Payload	PAYLOAD_MASS_KG	Orbit	Customer	Mission_Outcome
2010-04-06	18:45:00	F9 v1.0 B0003	CCAFS LC-40	Dragon Spacecraft Qualification Unit	0	LEO	SpaceX	Success
2010-08-12	15:43:00	F9 v1.0 B0004	CCAFS LC-40	Dragon demo flight C1, two CubeSats, barrel of Brouere cheese	0	LEO (ISS)	NASA (COTS) NRO	Success
2012-05-22	07:44:00	F9 v1.0 B0005	CCAFS LC-40	Dragon demo flight C2	525	LEO (ISS)	NASA (COTS)	Success
2012-08-10	00:35:00	F9 v1.0 B0006	CCAFS LC-40	SpaceX CRS-1	500	LEO (ISS)	NASA (CRS)	Success
2013-01-03	15:10:00	F9 v1.0 B0007	CCAFS LC-40	SpaceX CRS-2	677	LEO (ISS)	NASA (CRS)	Success

# Total Payload Mass

---

```
] : %sql select sum(payload_mass_kg_) from spacetable where customer='NASA (CRS)'  
* sqlite:///my_data1.db  
Done.  
]: sum(payload_mass_kg_)  
45596
```

# Average Payload Mass by F9 v1.1

---

```
%sql select avg(payload_mass_kg_) from spacetable where booster_version like 'F9 v1.1'
```

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

<b>avg(payload_mass_kg_)</b>
------------------------------

2534.6666666666665
--------------------

# First Successful Ground Landing Date

---

```
: %sql select min(date) from spacetable where landing_outcome='Success (ground pad)'  
* sqlite:///my_data1.db  
Done.  
: min(date)  
-----  
2015-12-22
```

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

---

```
}]: %sql select distinct(booster_version) from spacetable where landing_outcome='Success (drone ship)' and payload_mass_kg_ between 4000 and 6000
* sqlite:///my_data1.db
Done.
}]: 

| Booster_Version |
|-----------------|
| F9 FT B1022     |
| F9 FT B1026     |
| F9 FT B1021.2   |
| F9 FT B1031.2   |


```

# Total Number of Successful and Failure Mission Outcomes

---

## Task 7

List the total number of successful and failure mission outcomes

```
1]: %sql select mission_outcome, count(*) from spacetable group by mission_outcome
```

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

```
1]:
```

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



# Boosters Carried Maximum Payload

---

## Task 8

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

```
] : %sql select booster_version from spacetable where payload_mass_kg_ = (select max(payload_mass_kg_) from spacetable)
```

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

```
] : Booster_Version
```

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

# 2015 Launch Records

---

in year 2015:

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

```
%sql select substr(date, 6, 2) as month, booster_version, launch_site, landing_outcome from spacetable where substr(Date,0,5)='2015'
```

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

month	Booster_Version	Launch_Site	Landing_Outcome
10	F9 v1.1 B1012	CCAFS LC-40	Failure (drone ship)
04	F9 v1.1 B1015	CCAFS LC-40	Failure (drone ship)

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

---

in descending order.

```
%sql select landing_outcome, count(*) as cnt from spacetable group by landing_outcome having date between '2010-06-04' and '2017-03-20'
```

\* sqlite:///my\_data1.db  
Done.

Landing_Outcome	cnt
No attempt	21
Success (drone ship)	14
Success (ground pad)	9
Failure (drone ship)	5
Controlled (ocean)	5
Uncontrolled (ocean)	2
Precluded (drone ship)	1

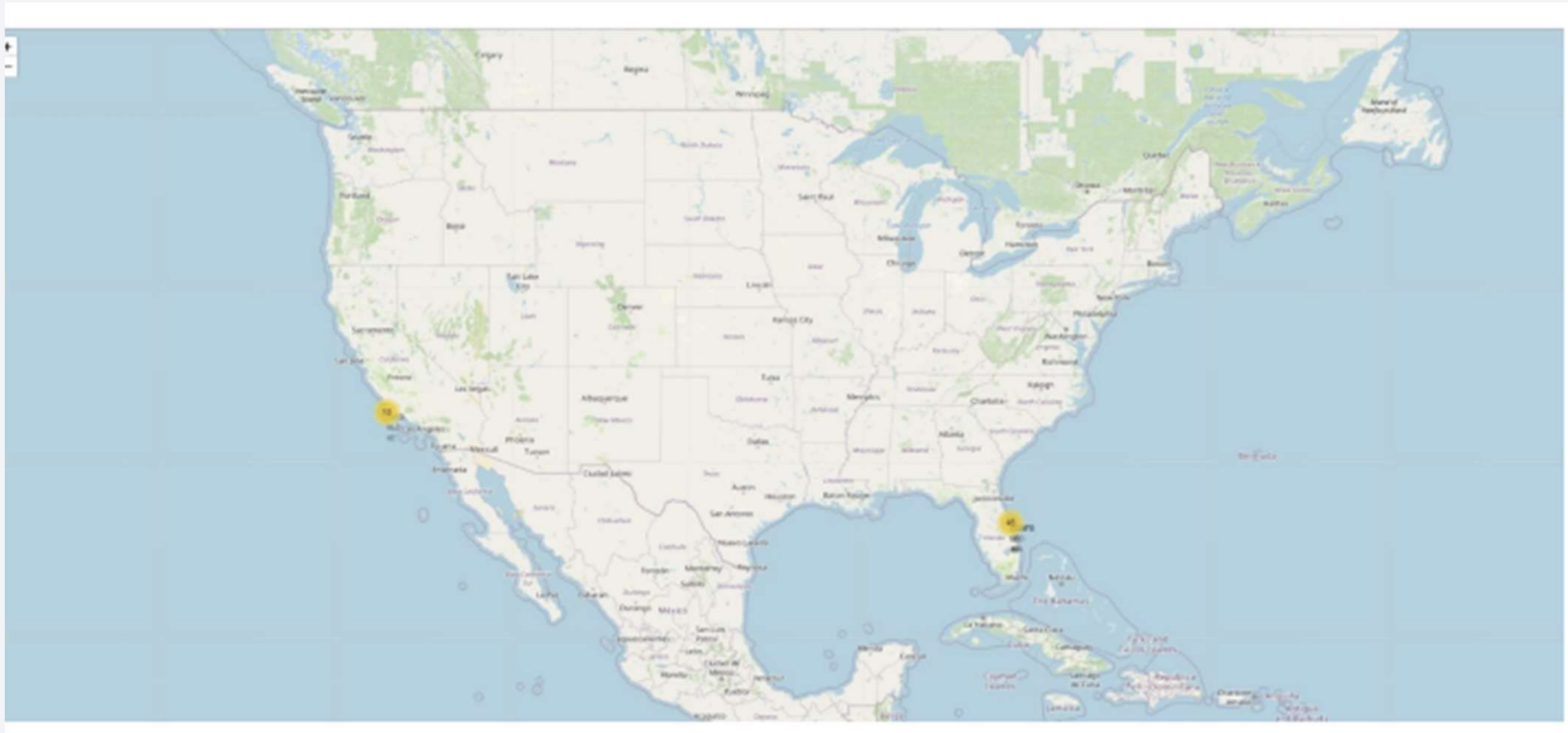
A satellite view of Earth from space, showing the curvature of the planet and the glow of city lights at night. The image is used as a background for the title slide.

Section 3

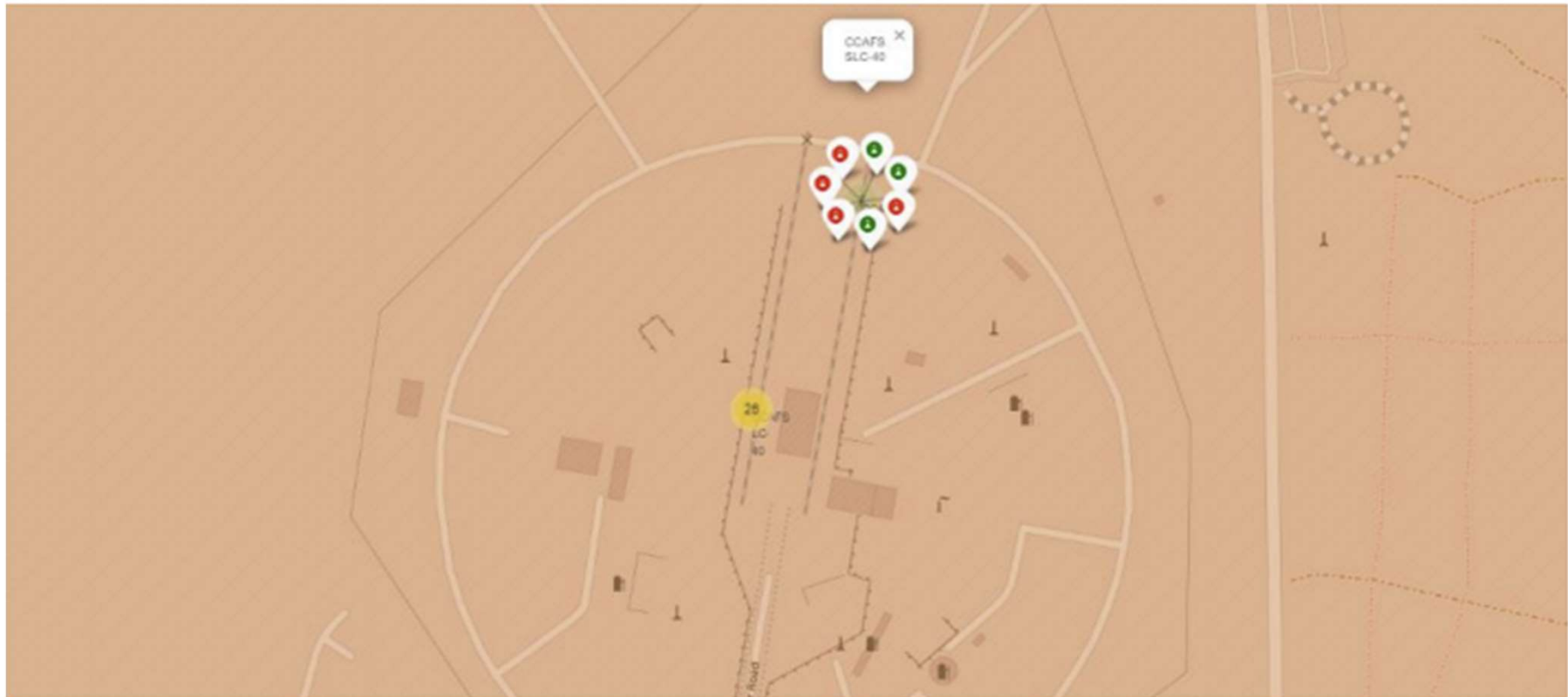
# Launch Sites Proximities Analysis

# Launch Sites

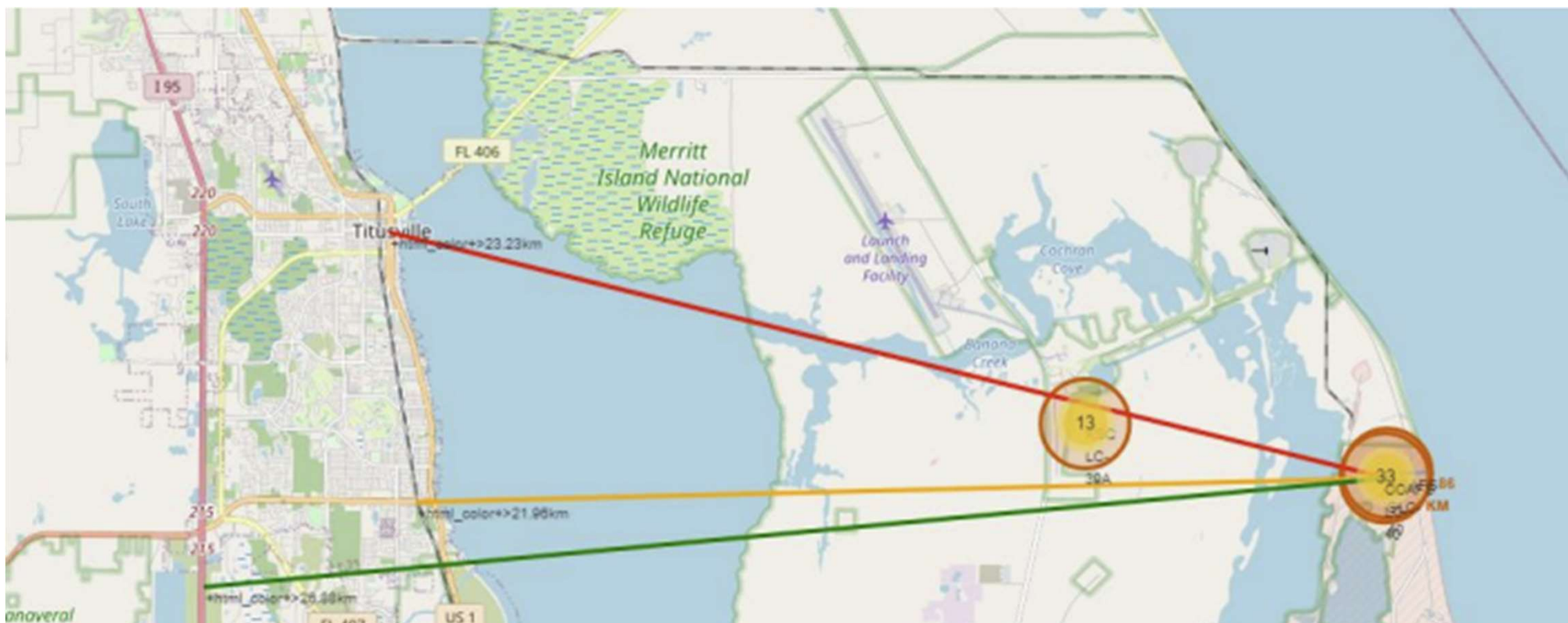
---



# Launch Outcomes



## Distance to Proximities







Section 4

# Build a Dashboard with Plotly Dash

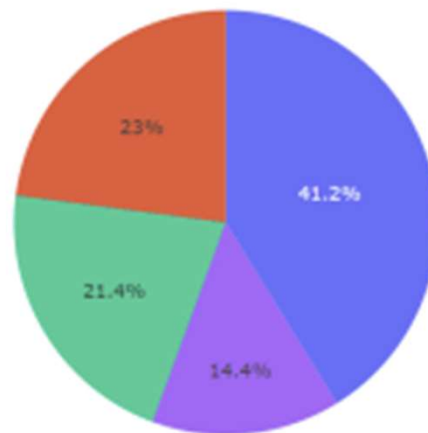
# Launch Success By Site

## SpaceX Launch Records Dashboard

All Sites



Total Success Launches by Site



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

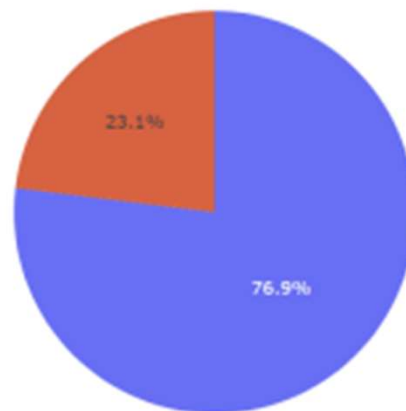
# Launch Success (KSC LC-29A)

## SpaceX Launch Records Dashboard

KSC LC-39A

× ▾

Total Success Launches for Site KSC LC-39A



■ 0  
■ 1

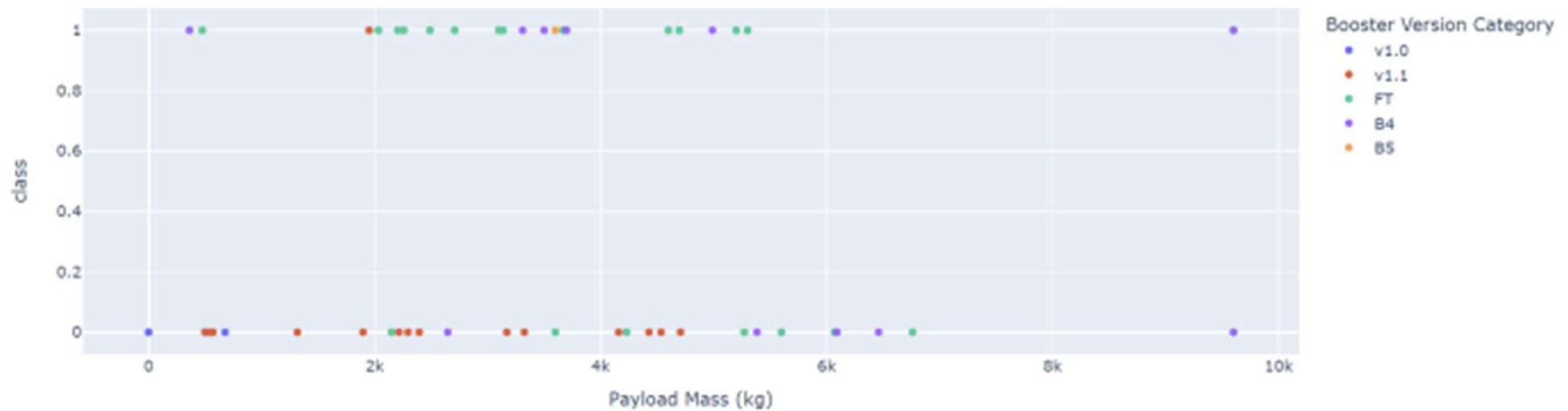
Class 0 = Fail  
Class 1 = Success

# Payload Mass and Success

Payload range (Kg):



Correlation Between Payload and Success for All Sites





Section 5

# Predictive Analysis (Classification)

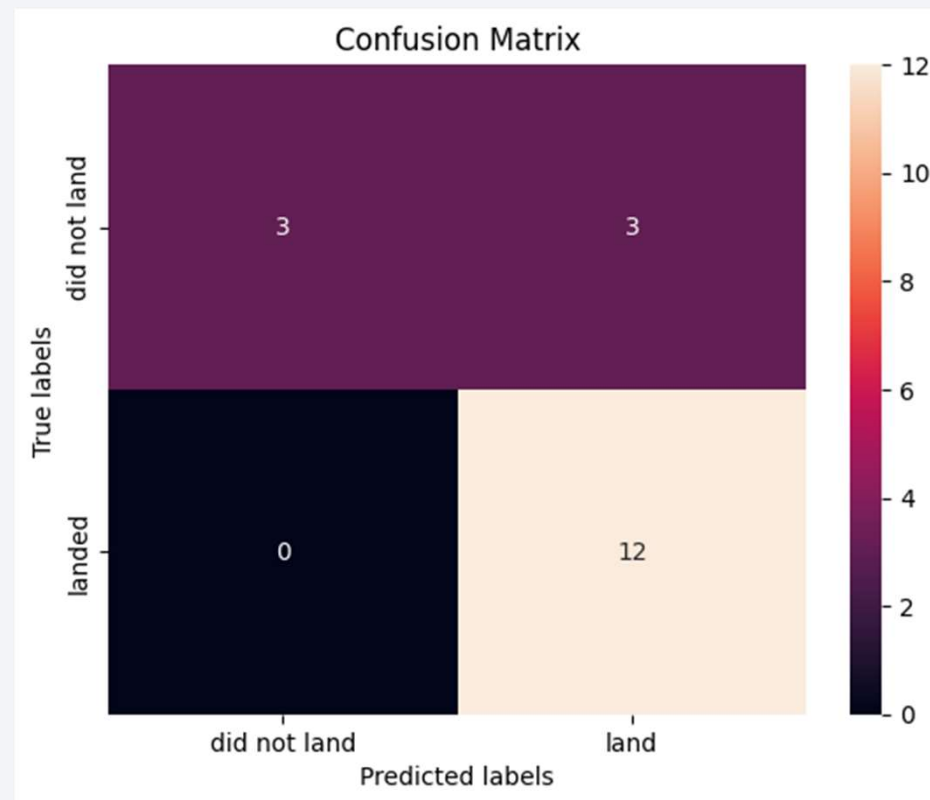
# Classification Accuracy

---

	LogReg	SVM	Tree	KNN
Jaccard_Score	0.800000	0.800000	0.800000	0.800000
F1_Score	0.888889	0.888889	0.888889	0.888889
Accuracy	0.833333	0.833333	0.833333	0.833333

# Confusion Matrix

---



# Conclusions

---

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



# Appendix

---

- [https://github.com/ScottLiao920/IBMDS\\_Coursera/tree/main](https://github.com/ScottLiao920/IBMDS_Coursera/tree/main)

Thank you!

