



IBM Developer
SKILLS NETWORK

Winning Space Race with Data Science

John Pauline Pineda
October 14, 2023

Outline

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

Executive Summary

- Methodologies explored in the study included the following:
 - Data collection using REST API and web scraping
 - Pre-processing methods
 - Exploratory data analysis using visualization and SQL
 - Interactive visual analytics using Folium and Plotly Dashboards
 - Predictive analysis encompassing model development, evaluation, and deployment
- Study results provided the following:
 - Complete and pre-processed data appropriate for analysis
 - Determined effects of features to successful landing outcomes
 - Determined effects of spatial features to successful launch outcomes
 - Determined effects of features to successful launch outcomes
 - Classification model for predicting landing outcomes with high accuracy

Introduction

- Space exploration has traditionally been associated with high costs. However, under the leadership of **Elon Musk**, has emerged as a cost-effective alternative.
- Notably, SpaceX has revolutionized the space exploration industry through their focus on **cost savings and cost-efficiency** that sets it apart from traditional aerospace companies.
- SpaceX pioneered the **reusable rocket technology**. Their ability to reuse the reusable first stage, allows multiple launches with the same hardware. By reusing portions of their launch vehicles, SpaceX dramatically reduces the overall mission expenses.
- This capstone project generally aims to investigate the factors influencing the success rate of SpaceX's first stage through exploratory data analysis and machine learning. By understanding these factors, thereby allowing reuse for multiple launches and achieving cost efficiency.
- In particular, a classification model will be formulated with the aim of predicting the outcome of the Falcon 9 rocket's successful landing based on various predictors, while delivering accurate predictions when a rocket lands successfully.

Section 1

Methodology

Methodology

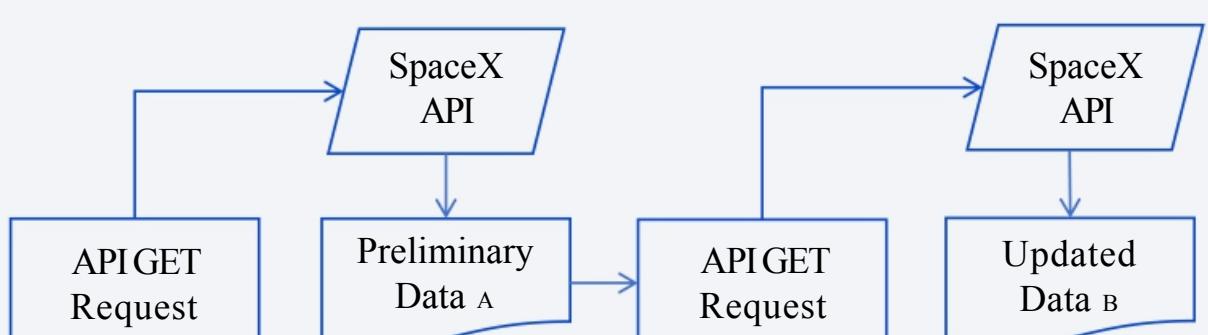
Executive Summary

- Data collection methodology
 - Sources: SpaceX REST API + web scraping
- Perform data wrangling
 - Pre-Processing: row + column filtering, missing data imputation
- Perform exploratory data analysis (EDA) using visualization tools
- Perform interactive visual analytics using Folium and Streamlit
- Perform predictive analysis using classification models
 - Classification Models: decision tree, k-nearest neighbors, logistic regression
 - Train-Test Ratio: 80% train and hyperparameter tuning
 - Model Performance Evaluation: accuracy on train data with cross-validation

Data Collection

- SpaceX launch data was gathered from 2 sources:
 - **Source 1:** SpaceX REST API containing data about Space rocket used, payload delivered, launch specifications, lan among others.
 - **Source 2:** web scraping of Falcon 9 launch data from rele
- Data collection process flowchart is presented as foll

Option 1: SpaceX REST API

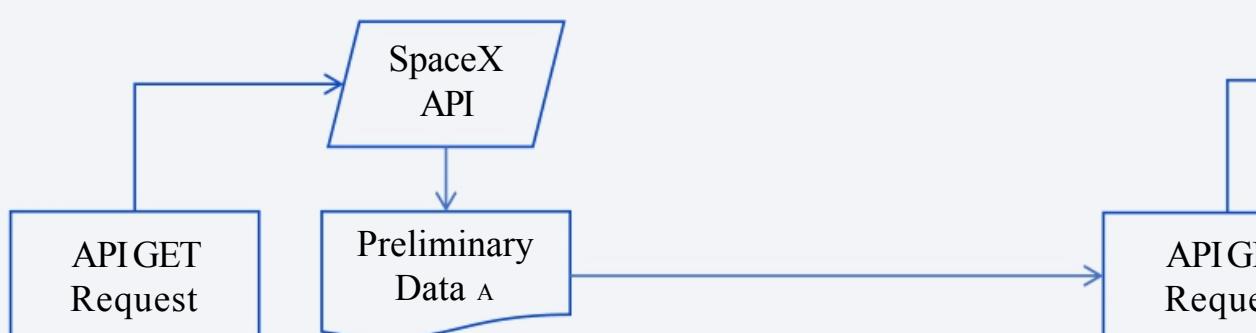


Data: A Rocket Launch B Rocket Launch + Booster + Launch Pad + Payload +

Data Collection – SpaceX API

- Data collection process flowchart with SpaceX REST

Option 1: SpaceX REST API



API GET Request Steps

- Define SpaceX URL for Rocket Launch data
- Apply API GET request on SpaceX URL
- Decode the API response as JSON object
- Transform JSON object to a Pandas data frame
- Perform necessary feature sub-setting and row filtering
- Generate preliminary data from filtered data frame
- Create additional empty global variable list

API GET

- Defi
+ Pa
- Defi
on th
- Impl
- Upd
- Upd

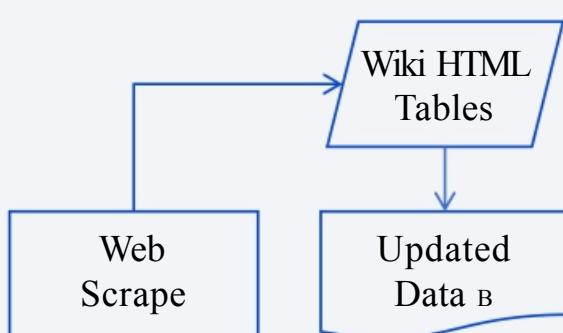
Data: _A Rocket Launch _B Rocket Launch + Booster + Launch Pad + Payload +

- GitHub URL of the completed SpaceX API calls Python

Data Collection – Scraping

- Data collection process flowchart with web scraping

Option 2: World Wide Web



Web Scrape Steps

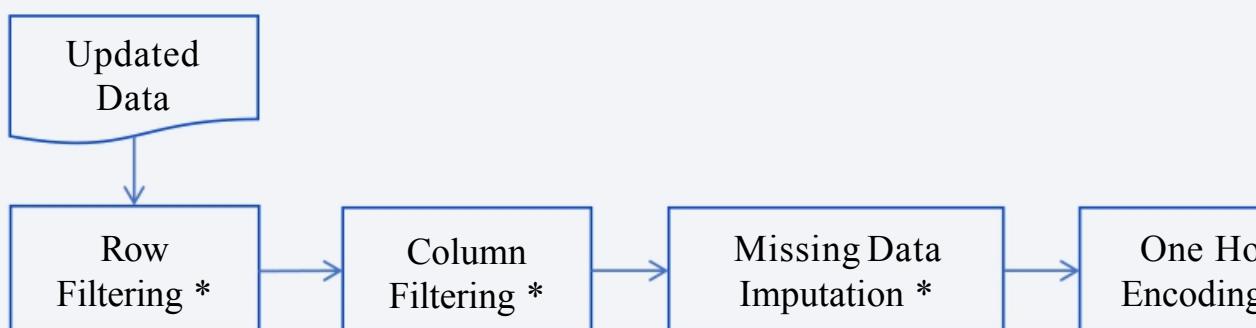
- Define WWW URL for Wiki HTML tables for Booster + Launch Pad + Payload + Core data
- Apply API GET request on WWW URL
- Decode the API response as BeautifulSoup object
- Extract all column names and row information from the BeautifulSoup object HTML tables
- Create a Pandas data frame from the parsed HTML tables

Data: B Rocket Launch + Booster + Launch Pad + Payload + Core

- GitHub URL of the completed web scraping Python n

Data Wrangling

- Data wrangling was applied on the updated data set and transforming raw data into a format suitable for machine learning.
- Data wrangling process flowchart is presented as follows:



* Data wrangling steps which were incorporated under the data collection Python notebook.

- GitHub URL of the data collection Python notebook will be provided [here](#).
- GitHub URL of the completed data wrangling Python notebook will be provided [here](#).

EDA with Data Visualization

- Exploratory data analysis involved formulating the analysis questions and visualize the effects of the different features on landing success rate
 - **Categorical Plot:** Effect of payload mass and flight number
 - **Categorical Plot:** Effect of launch site and flight number
 - **Categorical Plot:** Effect of payload mass and launch site
 - **Categorical Plot:** Effect of flight number and orbit to landing
 - **Categorical Plot:** Effect of payload mass and orbit to landing
 - **Bar Plot:** Effect of orbit to landing success rate
 - **Line Plot:** Effect of year to landing success rate
- GitHub URL of the completed EDA with data visualization

EDA with SQL

- Exploratory data analysis involved drill-down investigation
 - **Distinct Selection:** SELECT statement + DISTINCT expression
 - **String Patterns:** SELECT statement + WHERE clause + LIKE operator
 - **Functions:** SELECT statement + SUM|AVG |MIN functions
 - **Multiple Conditions:** SELECT statement + WHERE clause + AND|OR operators
 - **Grouping Result Sets:** SELECT statement + COUNT function
 - **Embedded Subquery:** SELECT statement + WHERE clause + IN operator
 - **Substring Extraction:** SELECT statement + SUBSTR function
 - **Range Conditions:** SELECT statement + COUNT function
 - **Sorting Result Sets:** SELECT statement + COUNT function
- GitHub URL of the completed EDA with SQL Python

Build an Interactive Map with Folium

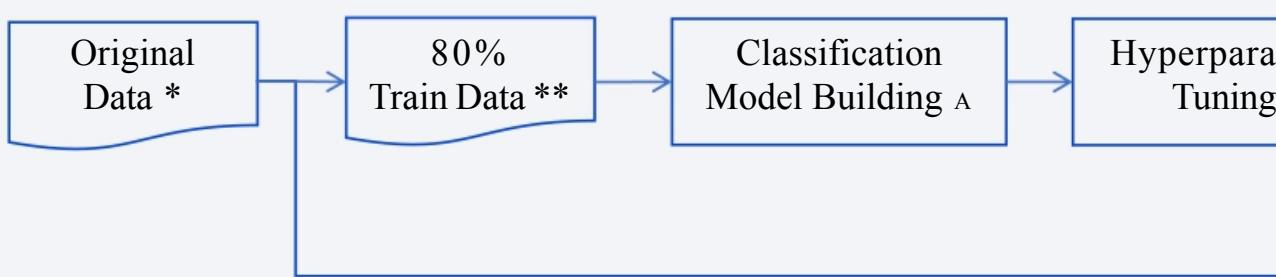
- Proximity analysis of launch data was conducted using Folium library
 - **Circles:** Indicates the launch sites
 - **Markers:** Indicates the names of the launch sites as a pop-up
 - **Marker Cluster:** Indicates individual launch record for each site
 - **Mouse Position:** Indicates latitude and longitude information
 - **Lines:** Connects the launch sites to identified landmarks (e.g., Cape Canaveral)
 - **Distance Circle:** Indicates the distance between the launch sites
- GitHub URL of the completed interactive Folium map
- **NOTE:** Despite enabling the ‘Trust Notebook’ setting, the uploaded python notebook is viewed through GitHub. Thus, the map outputs were provided [here](#) as an alternative.

Build a Dashboard with Plotly

- Interactive visual analytics involved formulating the research question and hypotheses to investigate the individual and combined effects of different variables.
- **Pie Chart:** Contribution of each launch site on the combined success rate.
- **Pie Chart:** Distribution of the launch outcomes for each launch site.
- **Categorical Plot:** Effect of payload mass and booster version on success rate.
- **Range Slider:** Effect of changing the payload mass range on success rate.
- GitHub URL of the completed Plotly Dash Python code.

Predictive Analysis (Classification)

- Predictive analysis involved the process steps described below:
 - **Model Development:** Included classification model development
 - **Model Evaluation:** Involved both internal and external evaluation
 - **Model Selection:** Involved the identification of the best performing model

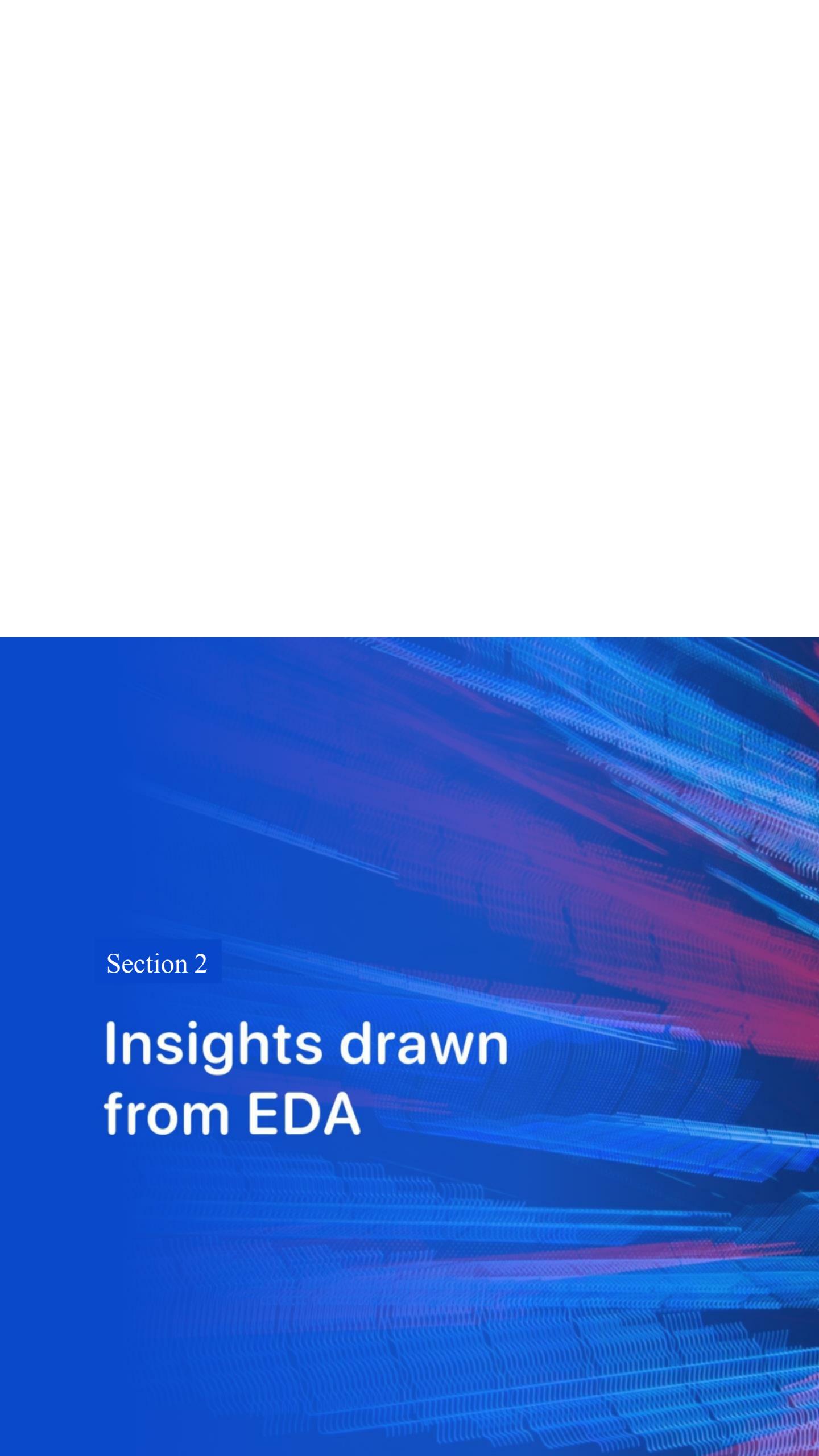


* 90 observations + 83 features | ** 72 observations + 83 features | *** 18 observations
^ 4 candidate models built using the Logistic Regression, Support Vector Machine, Random Forest, and Naive Bayes classifiers
^ Hyperparameters of 4 candidate models fine-tuned and internally evaluated using 5-fold cross-validation
C 4 candidate models with optimal hyperparameters were externally evaluated using 10-fold cross-validation
D Best model among 4 optimal candidate models determined based on internal evaluation metrics

- GitHub URL of the completed Predictive analysis Python code

Results

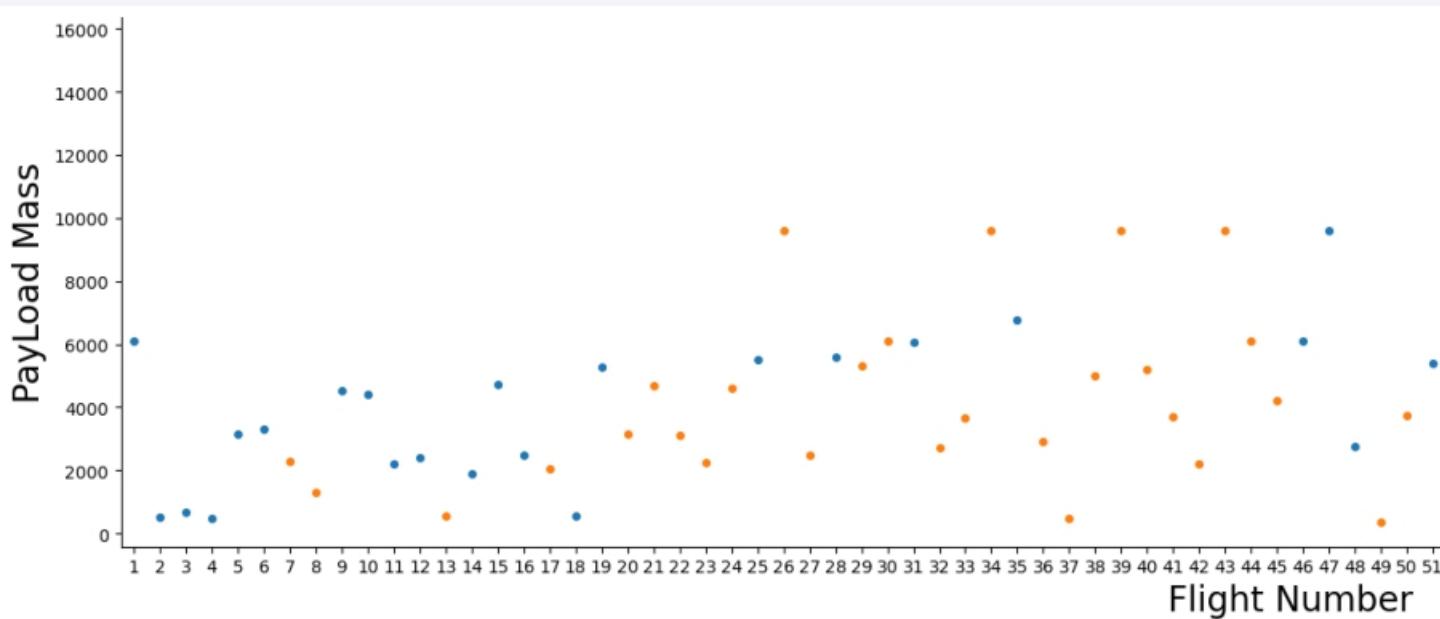
- Exploratory data analysis results
 - [Insights Drawn from EDA](#)
- Interactive analytics demo in screenshots
 - [Launch Sites Proximities Analysis](#)
 - [Build a Dashboard with Plotly Dash](#)
- Predictive analysis results
 - [Predictive Analysis \(Classification\)](#)

The background of the slide features a dynamic, abstract pattern of wavy, horizontal lines in shades of blue and red. These lines create a sense of depth and motion, resembling a digital or architectural landscape. They are more concentrated on the right side of the slide.

Section 2

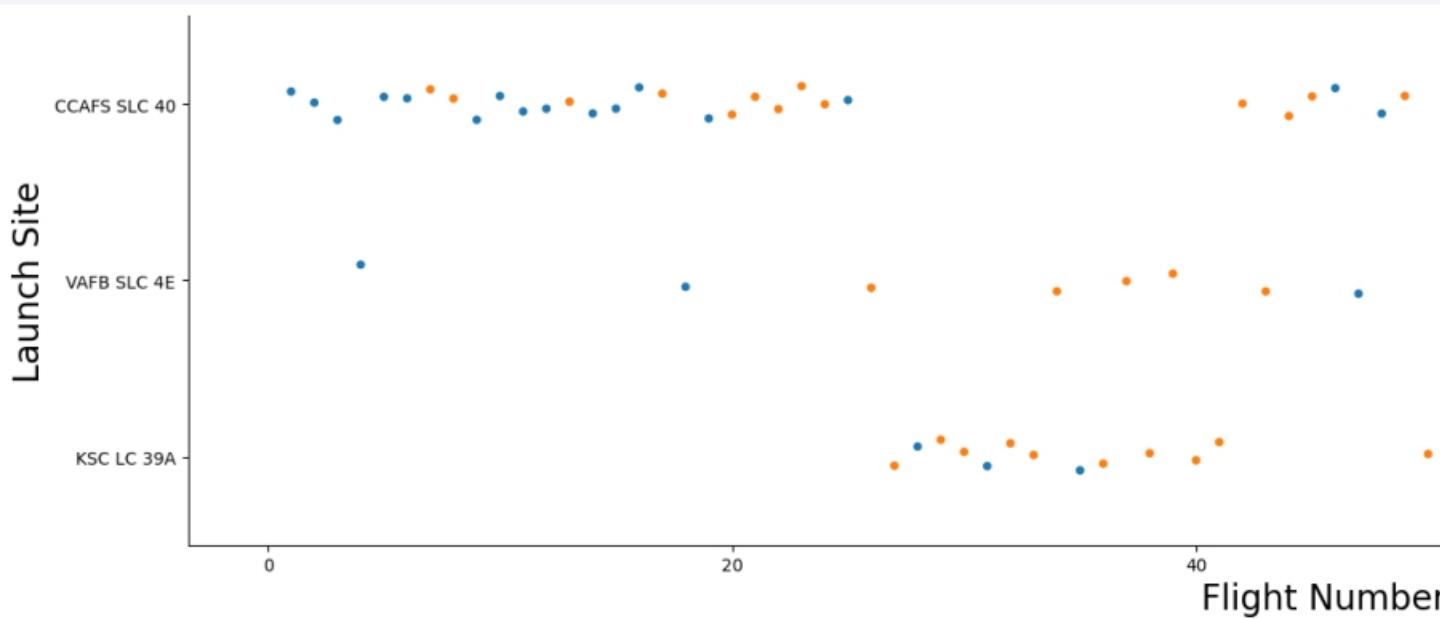
Insights drawn from EDA

Flight Number vs. Payload



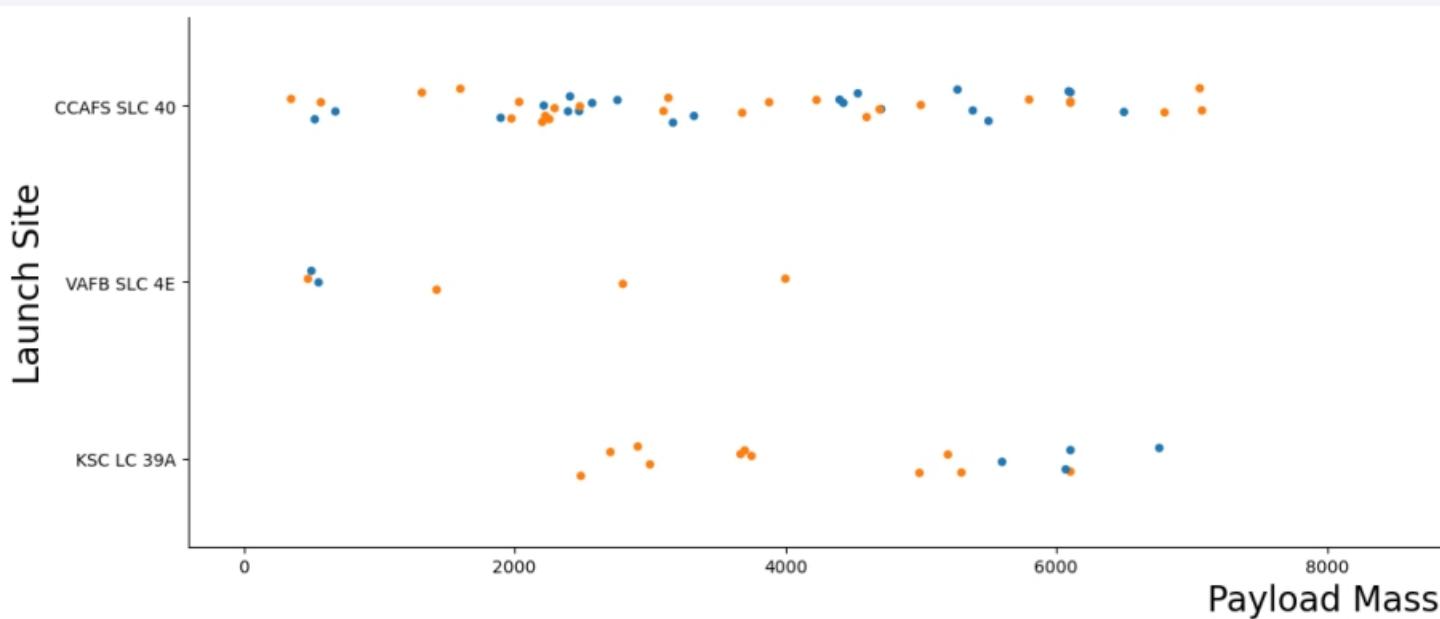
- Categorical Plot: Effect of Payload Mass and Flight N
 - More successful landings generally observed as flight nu
 - Successful landings observed across varying ranges of p

Flight Number vs. Launch Site



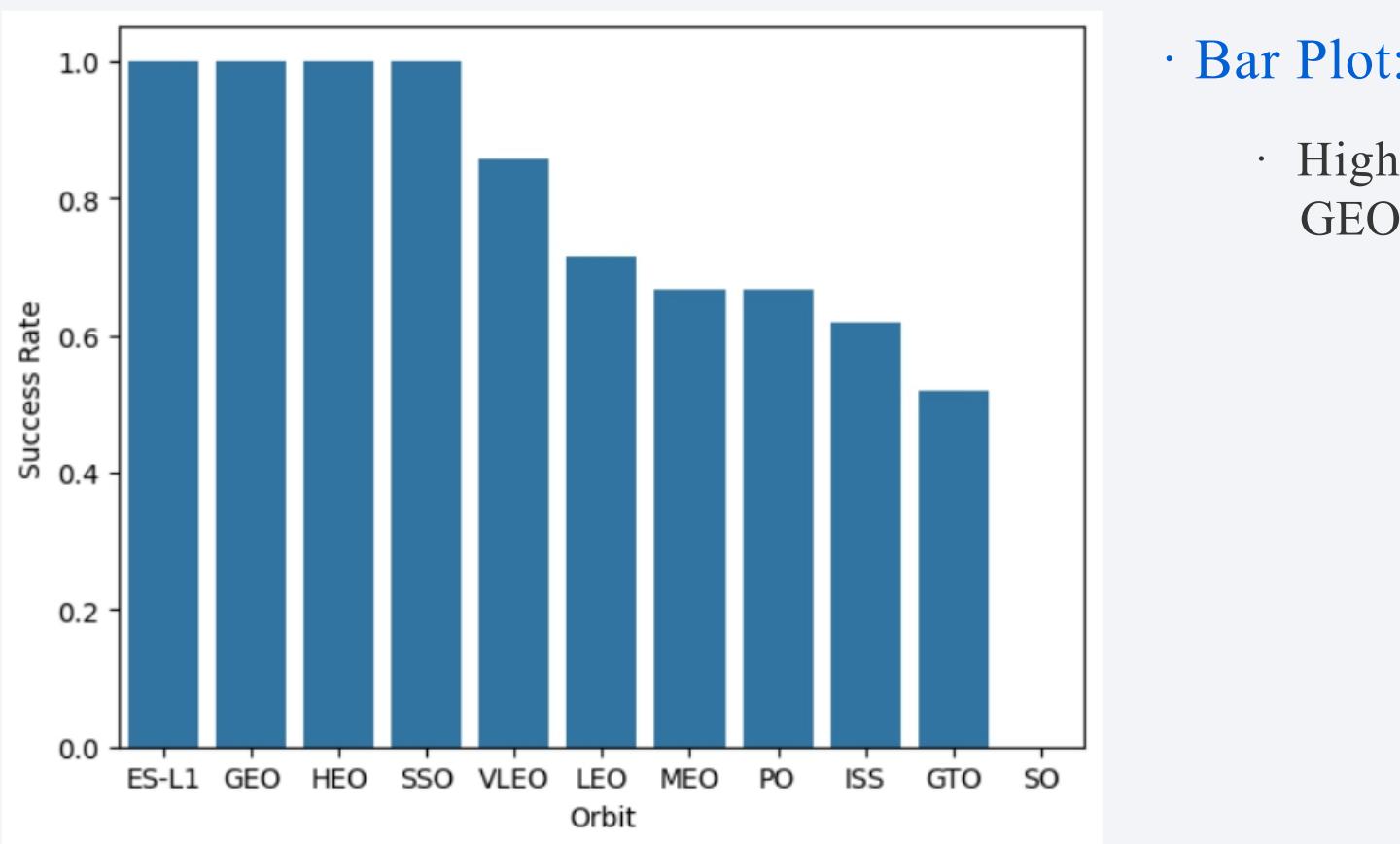
- **Categorical Plot:** Effect of Launch Site and Flight Number
- More successful landings generally observed for VAFB SLC 4E
- Increased flight numbers only observed for CCAFS SLC 40

Payload vs. Launch Site

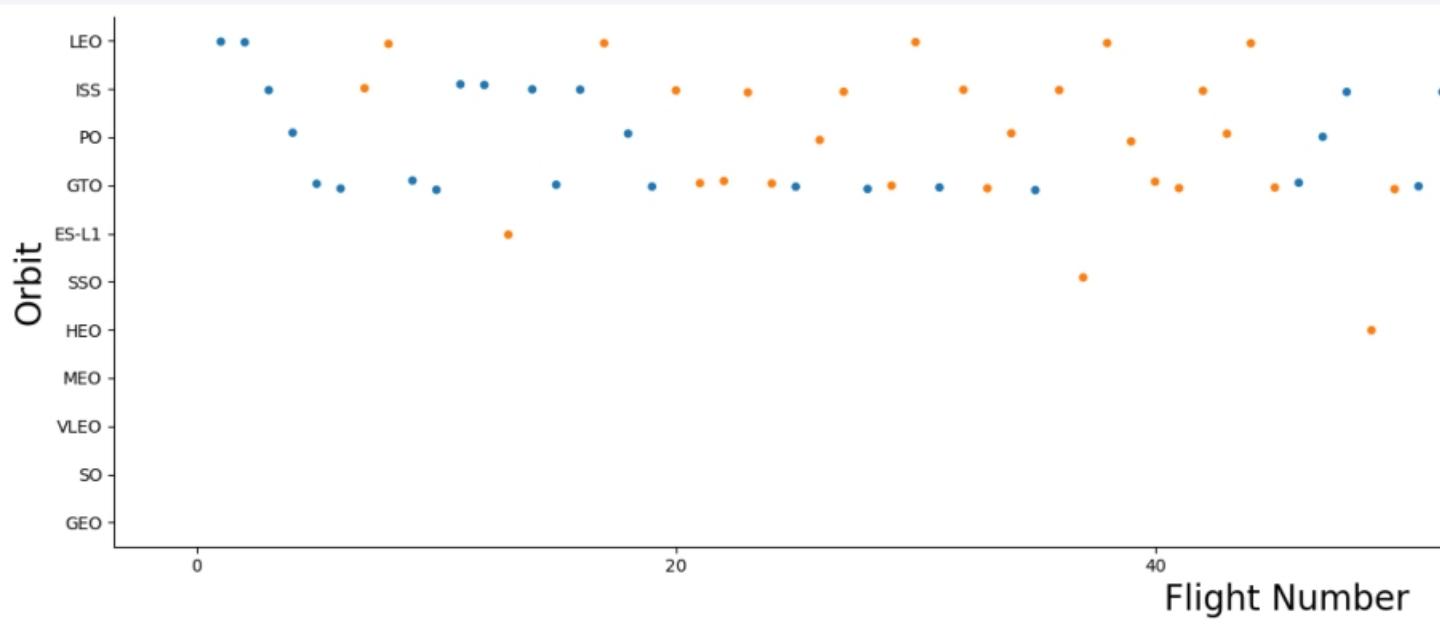


- **Categorical Plot:** Effect of Payload Mass and Launch
- Successful landings observed across varying ranges of p
- Higher payload masses >10K only observed for CCAFS S

Success Rate vs. Orbit Type

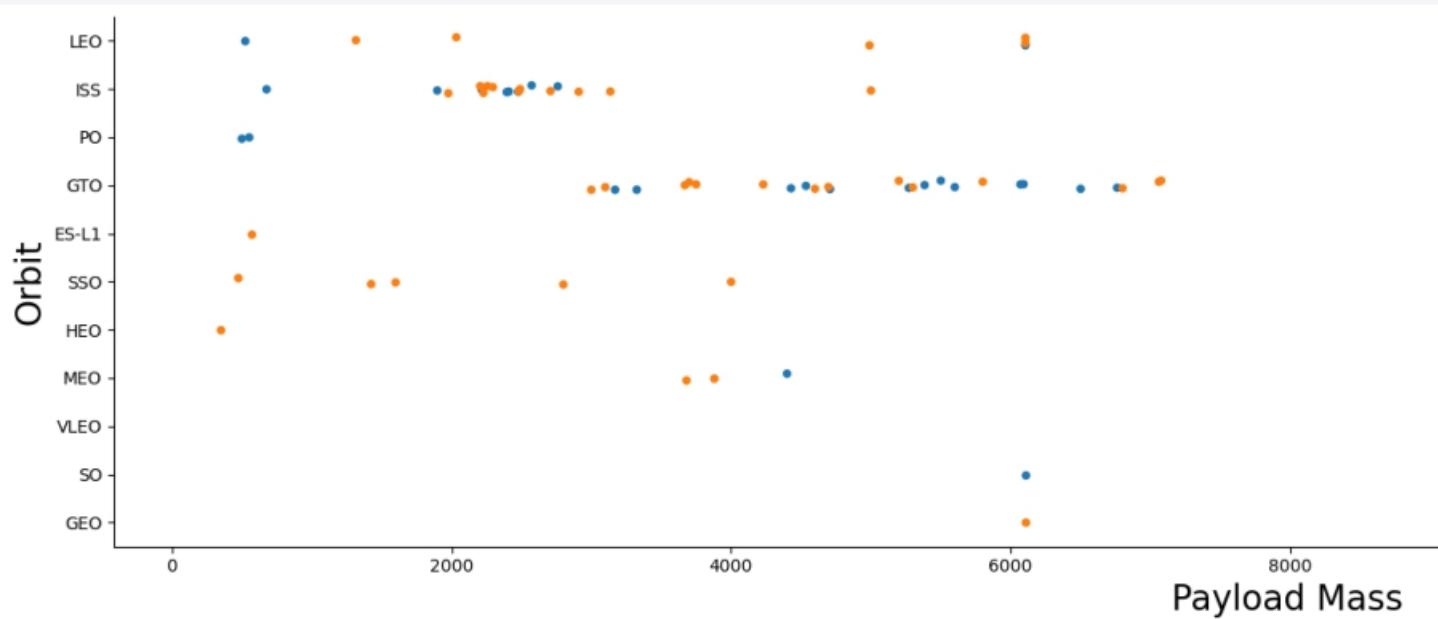


Flight Number vs. Orbit Type



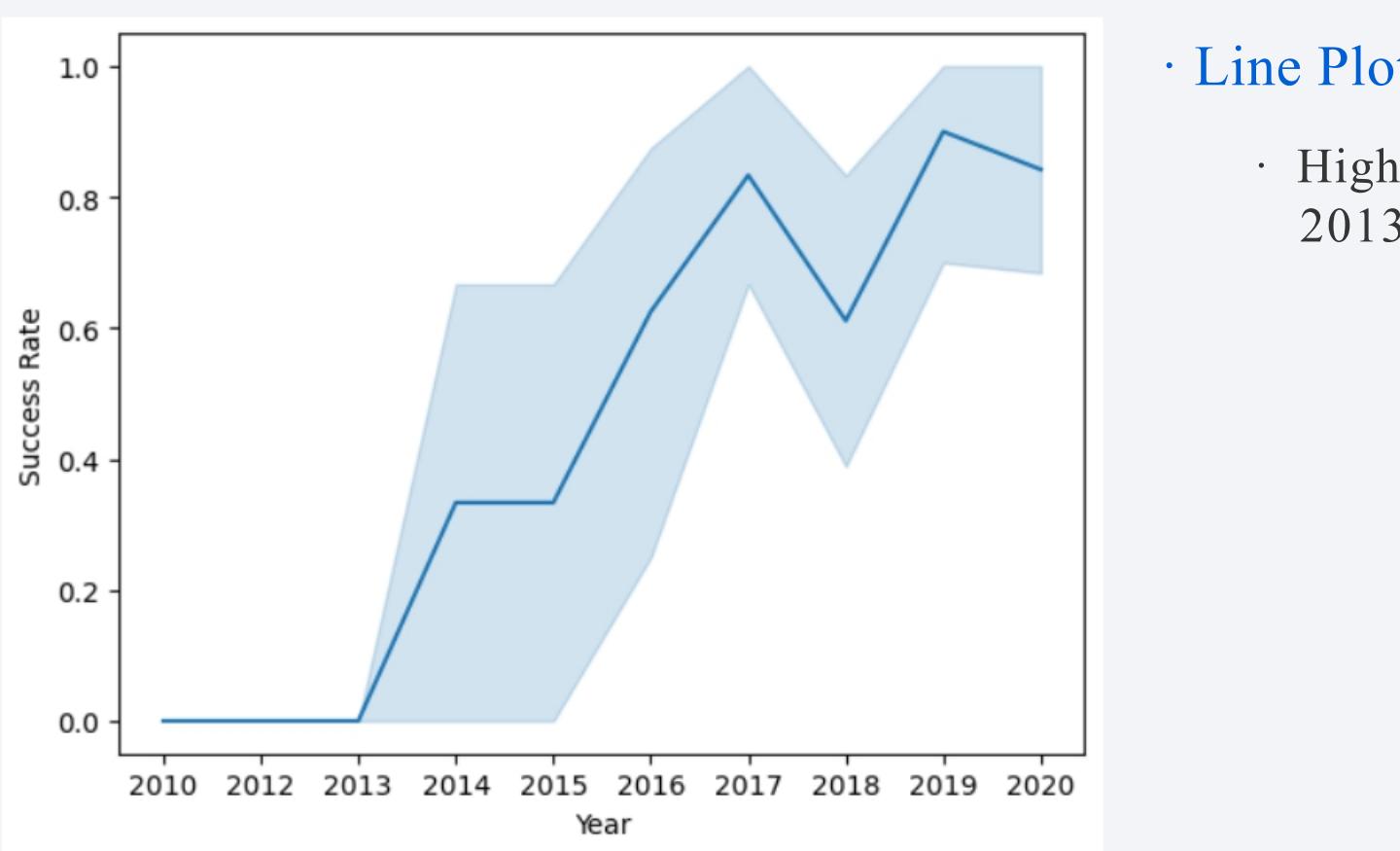
- **Categorical Plot:** Effect of Flight Number and Orbit type
 - Previous observation on higher landing success rate for E
 - Increased flight numbers only observed for the VLEO orb
 - Given a considerable number of flights, LEO, SSO and V

Payload vs. Orbit Type



- Categorical Plot: Effect of Payload Mass and Orbit type
 - Higher payload masses >9K only observed for ISS, PO and GTO orbits
 - Successful landings observed across varying ranges of payload mass

Launch Success Yearly Trend



All Launch Site Names

```
%%sql SELECT DISTINCT Launch_Site  
FROM SPACEXTABLE;
```

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

```
Done.
```

Launch_Site

CCAFS LC-40

VAFB SLC-4E

KSC LC-39A

CCAFS SLC-40

- **Distinct Selection:** SELECT statement + DISTINCT expression
 - There were 4 unique launch sites included the space mission. KSC LC-39A and CCAFS SLC-40.

Launch Site Names Begin with

```
%%sql
SELECT * FROM SPACEXTABLE
WHERE Launch_Site LIKE 'CCA%' LIMIT 5;
```

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

```
Done.
```

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

- **String Patterns:** SELECT statement + WHERE clause +
 - 5 records of launch sites beginning with the string ‘CCA’

Total Payload Mass

```
%%sql
SELECT SUM(PAYLOAD_MASS__KG_)
AS "TOTAL PAYLOAD MASS", Customer
FROM SPACEXTABLE
WHERE Customer LIKE "%NASA (CRS)%";
```

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

```
Done.
```

TOTAL PAYLOAD MASS	Customer
48213	NASA (CRS)

- **Functions:** SELECT statement + SUM function + WHERE clause
 - The total payload mass carried by boosters launched by NASA (CRS)

Average Payload Mass by F9 v1.1

```
%%sql
SELECT AVG(PAYLOAD_MASS__KG_)
AS "AVERAGE PAYLOAD MASS", Customer, Booster_Version
FROM SPACEXTABLE
WHERE Booster_Version LIKE 'F9 v1.1%';

* sqlite:///my_data1.db
Done.

AVERAGE PAYLOAD MASS  Customer  Booster_Version
2534.6666666666665      MDA      F9 v1.1 B1003
```

- **Functions:** SELECT statement + AVG function + WHERE clause
 - The average payload mass carried by booster version F9 v1.1

First Successful Ground Landi

```
%%sql
SELECT MIN(DATE) FROM SPACEXTABLE
WHERE "Landing_Outcome"
LIKE "Success (ground pad)"

* sqlite:///my_data1.db
Done.

MIN(DATE)
2015-12-22
```

- **Functions:** SELECT statement + MIN function + WHERE clause
 - The first successful landing outcome in a ground pad was on December 22, 2015.

Successful Drone Ship Landing between 4000 and 6000

```
%%sql
```

```
SELECT * FROM SPACEXTABLE  
WHERE Landing_Outcome = "Success (drone ship)" AND PAYLOAD_MASS_KG > 4000 AND PAYLOAD_MASS_KG < 6000
```

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

```
Done.
```

Date	Time (UTC)	Booster_Version	Launch_Site	Payload	PAYLOAD_MASS_KG	Orbit
2016-06-05	05:21:00	F9 FT B1022	CCAFS LC-40	JCSAT-14	4696	GTO
2016-08-14	05:26:00	F9 FT B1026	CCAFS LC-40	JCSAT-16	4600	GTO
2017-03-30	22:27:00	F9 FT B1021.2	KSC LC-39A	SES-10	5300	GTO
2017-11-10	22:53:00	F9 FT B1031.2	KSC LC-39A	SES-11 / EchoStar 105	5200	GTO

- **Multiple Conditions:** SELECT statement + WHERE clause
 - There were 4 boosters launched on the dates June 5, 2016, August 14, 2016, March 30, 2017, November 10, 2017 which had success landing in drone ship category and payload mass between 4000 and 6000.

Total Number of Successful and Failure Mission Outcomes

```
%%sql
SELECT Mission_Outcome, COUNT(Mission_Outcome) AS "Total Number of Mission Outcomes"
FROM SPACEXTABLE
GROUP BY Mission_Outcome;

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

Mission_Outcome	Total Number of Mission Outcomes
Failure (in flight)	1
Success	98
Success	1
Success (payload status unclear)	1

- **Grouping Result Sets:** SELECT statement + COUNT function
 - There were 100 successful mission outcomes and 1 failed mission outcome.

Boosters Carried Maximum Payload

```
%%sql
SELECT Booster_Version, Payload, PAYLOAD_MASS__KG_
FROM SPACEXTABLE WHERE PAYLOAD_MASS__KG_ = (SELECT MAX(PAYLOAD_MASS__KG_) FROM SPACEXTABLE);
* sqlite:///my_data1.db
Done.
```

Booster_Version	Payload	PAYLOAD_MASS__KG_
F9 B5 B1048.4	Starlink 1 v1.0, SpaceX CRS-19	15600
F9 B5 B1049.4	Starlink 2 v1.0, Crew Dragon in-flight abort test	15600
F9 B5 B1051.3	Starlink 3 v1.0, Starlink 4 v1.0	15600
F9 B5 B1056.4	Starlink 4 v1.0, SpaceX CRS-20	15600
F9 B5 B1048.5	Starlink 5 v1.0, Starlink 6 v1.0	15600
F9 B5 B1051.4	Starlink 6 v1.0, Crew Dragon Demo-2	15600
F9 B5 B1049.5	Starlink 7 v1.0, Starlink 8 v1.0	15600
F9 B5 B1060.2	Starlink 11 v1.0, Starlink 12 v1.0	15600
F9 B5 B1058.3	Starlink 12 v1.0, Starlink 13 v1.0	15600
F9 B5 B1051.6	Starlink 13 v1.0, Starlink 14 v1.0	15600
F9 B5 B1060.3	Starlink 14 v1.0, GPS III-04	15600
F9 B5 B1049.7	Starlink 15 v1.0, SpaceX CRS-21	15600

• Embedded WHERE

- The boosters carried maximum payload
- 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

```
%%sql
SELECT SUBSTR(Date, 6, 2) AS MONTH, Booster_Version, Launch_Site, Landing_Outcome
FROM SPACEXTABLE
WHERE SUBSTR(Date, 0, 5) = "2015" AND Landing_Outcome = 'Failure (drone ship)';
* 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)
```

- **Substring Extraction:** SELECT statement + SUBSTR function
 - There were two records for the months of April and October, month names, failure landing outcomes in drone ship, b

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

```
%%sql
SELECT Landing_Outcome, COUNT(Landing_Outcome) AS "Landing Outcome Count"
FROM SPACEXTABLE WHERE Date BETWEEN '2010-06-04' AND '2017-03-20'
GROUP BY Landing_Outcome
ORDER BY "Landing Outcome Count" DESC;
```

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

Landing_Outcome	Landing Outcome Count
No attempt	10
Success (ground pad)	5
Success (drone ship)	5
Failure (drone ship)	5
Controlled (ocean)	3
Uncontrolled (ocean)	2
Precluded (drone ship)	1
Failure (parachute)	1

- Range Condition
SELECT ... WHERE ...
clause +
BY | OR

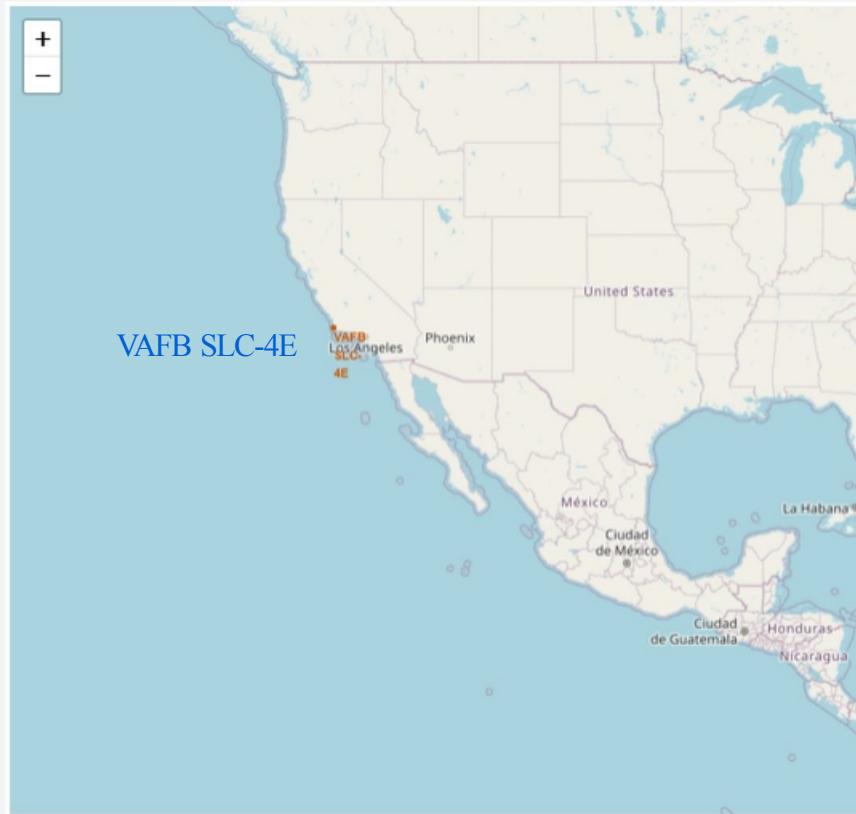
- There are many ways to do this
the details will depend on which

The background of the slide is a blue-toned satellite image of Earth's surface, showing landmasses and city lights at night.

Section 3

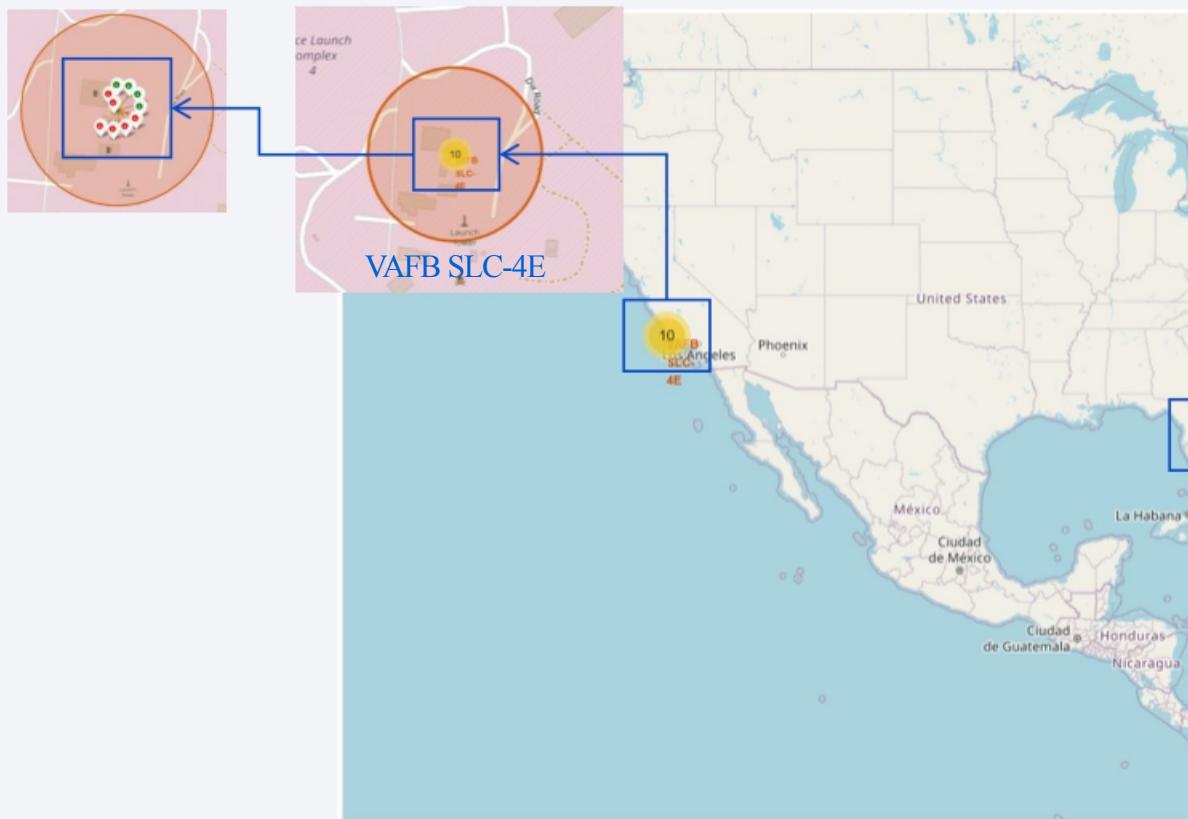
Launch Sites Proximities Analysis

Location Markers of Launch S



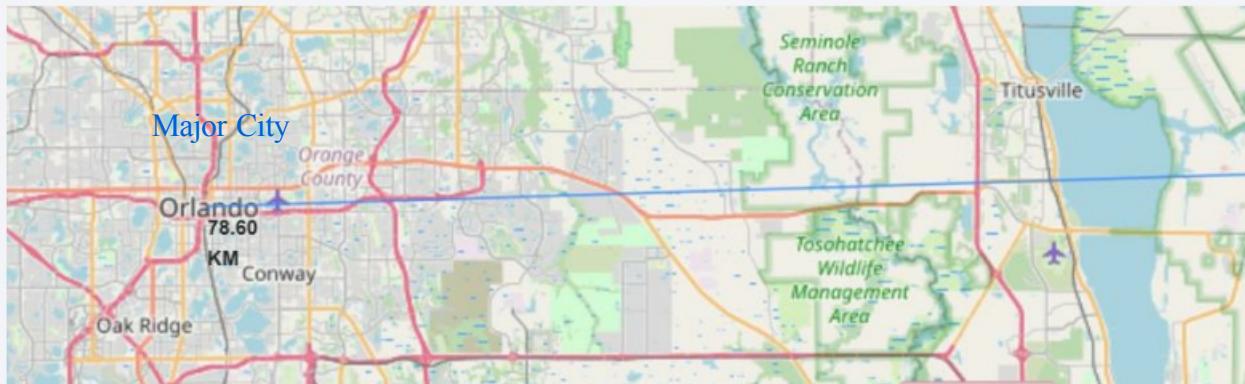
- **Location Markers:** US launch sites were normally situated
 - West Coast: VAFB SLC-4E
 - East Coast: KSC LC-39A + CCAFS SLC-40 + CCAFS LC-41

Color-Labeled Launch Outcomes

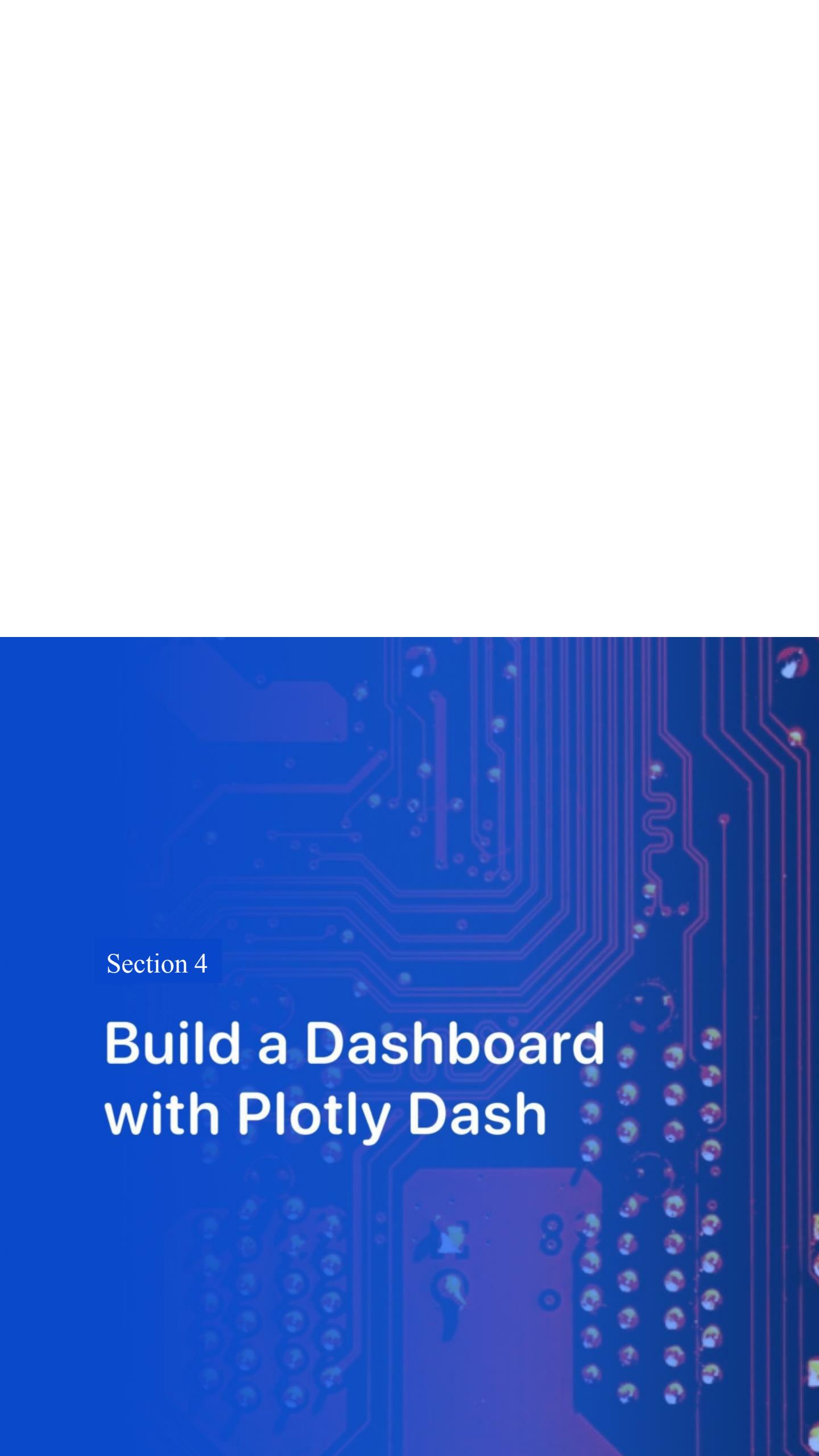


- **Outcome Markers:** Launch success varies per site.
- The KSC LC-39A site demonstrated a relatively higher rate of success compared to the CCAFS SLC-40, CCAFS LC-40 and VAFB SLC-4E sites.

Proximity Analysis of Launch Site Coastline, Highway, Railway and Waterways



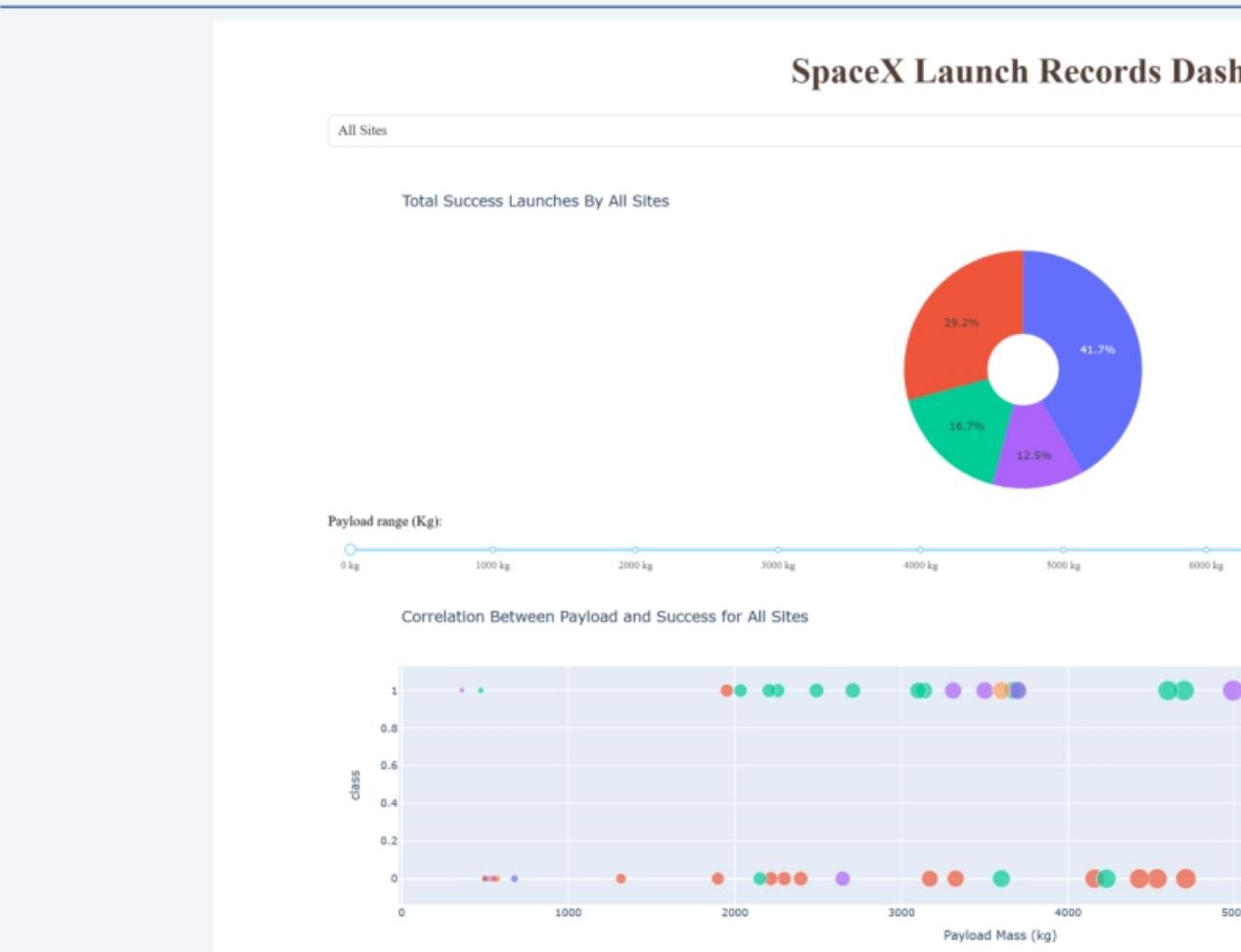
- **Distance Markers:** Launch sites were located far from coastlines, highways and railroads.
 - Taking CCAFS SLC-40 as an example, its approximate distance from Orlando was 78.60 Km.
 - The approximate distance of the CCAFS SLC-40 launch site from the coastline was 100 Km.
 - The approximate distance of the CCAFS SLC-40 launch site from the railway line was 100 Km.
 - The approximate distance of the CCAFS SLC-40 launch site from the highway was 100 Km.



Section 4

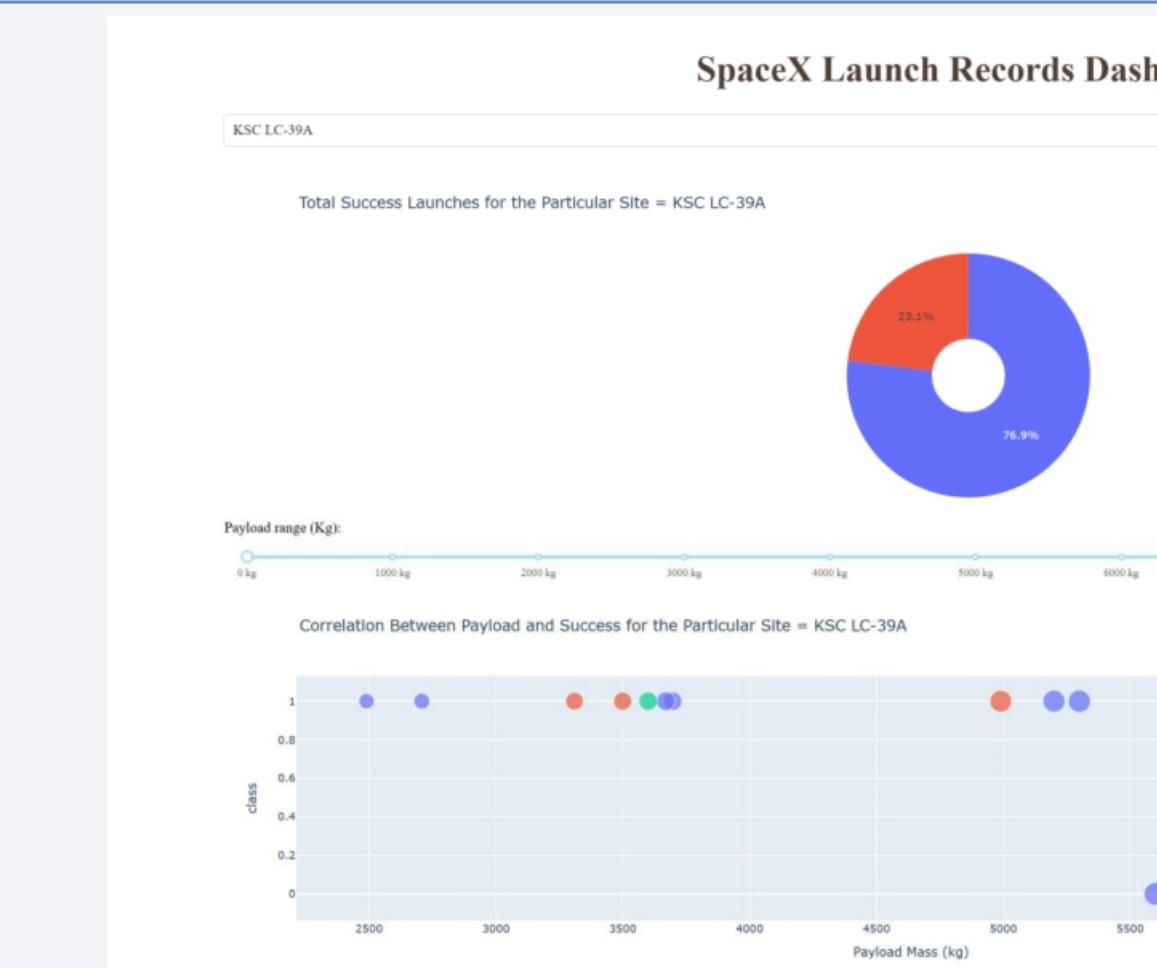
Build a Dashboard with Plotly Dash

Combined Launch Success Count



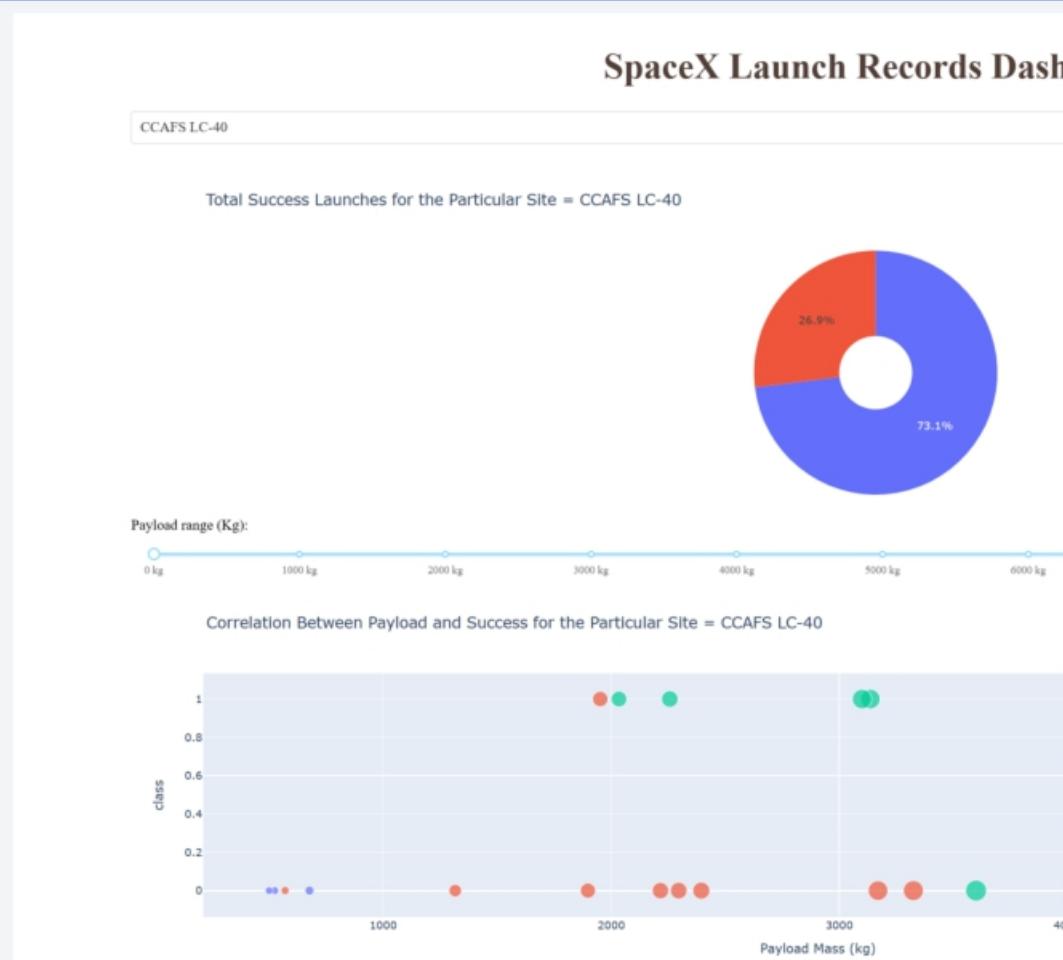
- **Launch Success Count:** The most number of successful launches were from KSC LC-39A contributing 41.7% to the overall launch count.

Launch Success Ratio Per Indi



- **Launch Success Ratio:** Among all sites, the launch success ratio for KSC LC-39A is ranked first. 76.9% of all launches from KSC LC-39A were successful.

Launch Success Ratio Per Indi



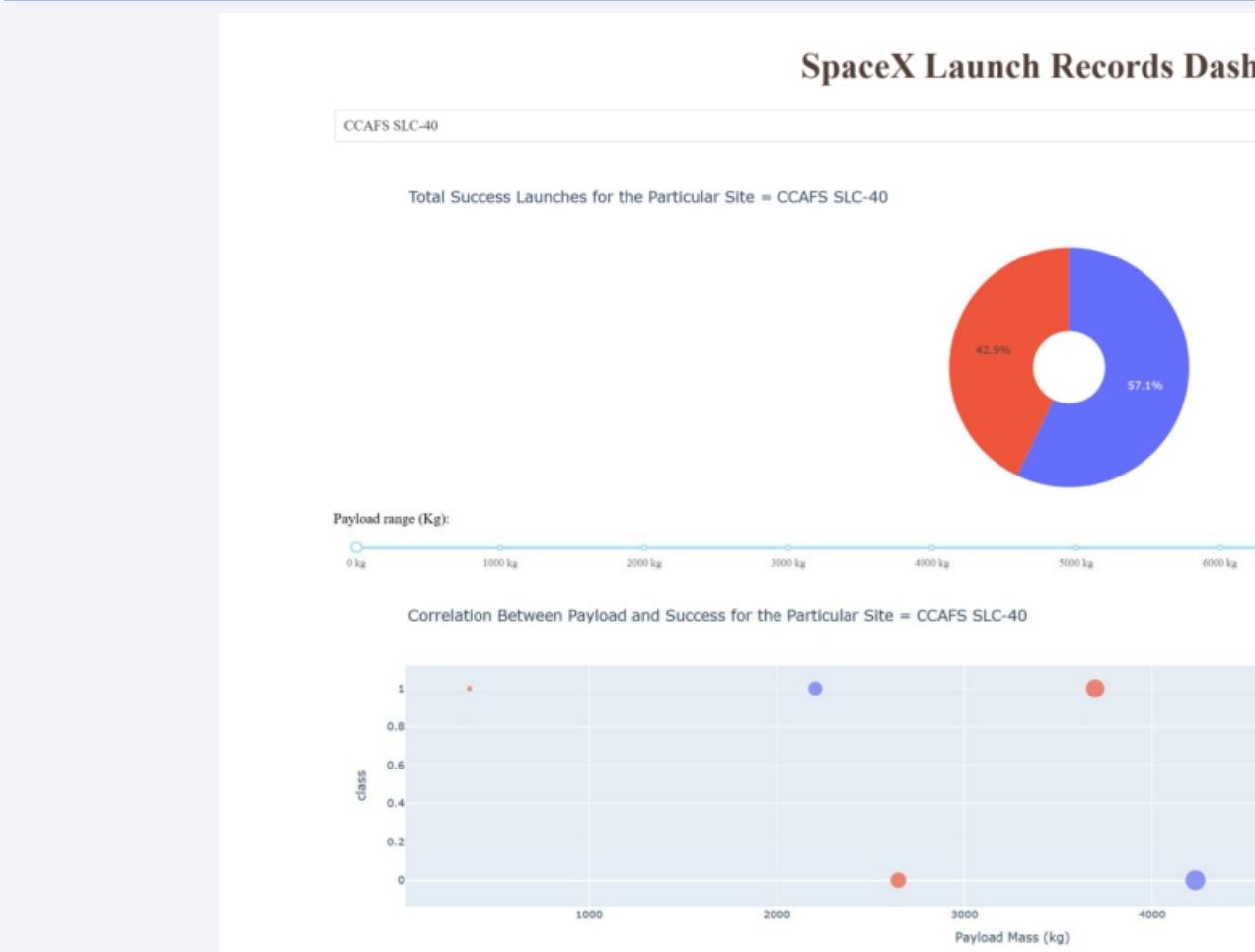
- **Launch Success Ratio:** Among all sites, the launch success ratio for CCAFS LC-40 is ranked second. 73.1% of all launches from CCAFS LC-40 were successful.

Launch Success Ratio Per Indi



- **Launch Success Ratio:** Among all sites, the launch success ratio for VAFB SLC-4E was ranked third. 60.0% of all launches from VAFB SLC-4E were successful.

Launch Success Ratio Per Indi



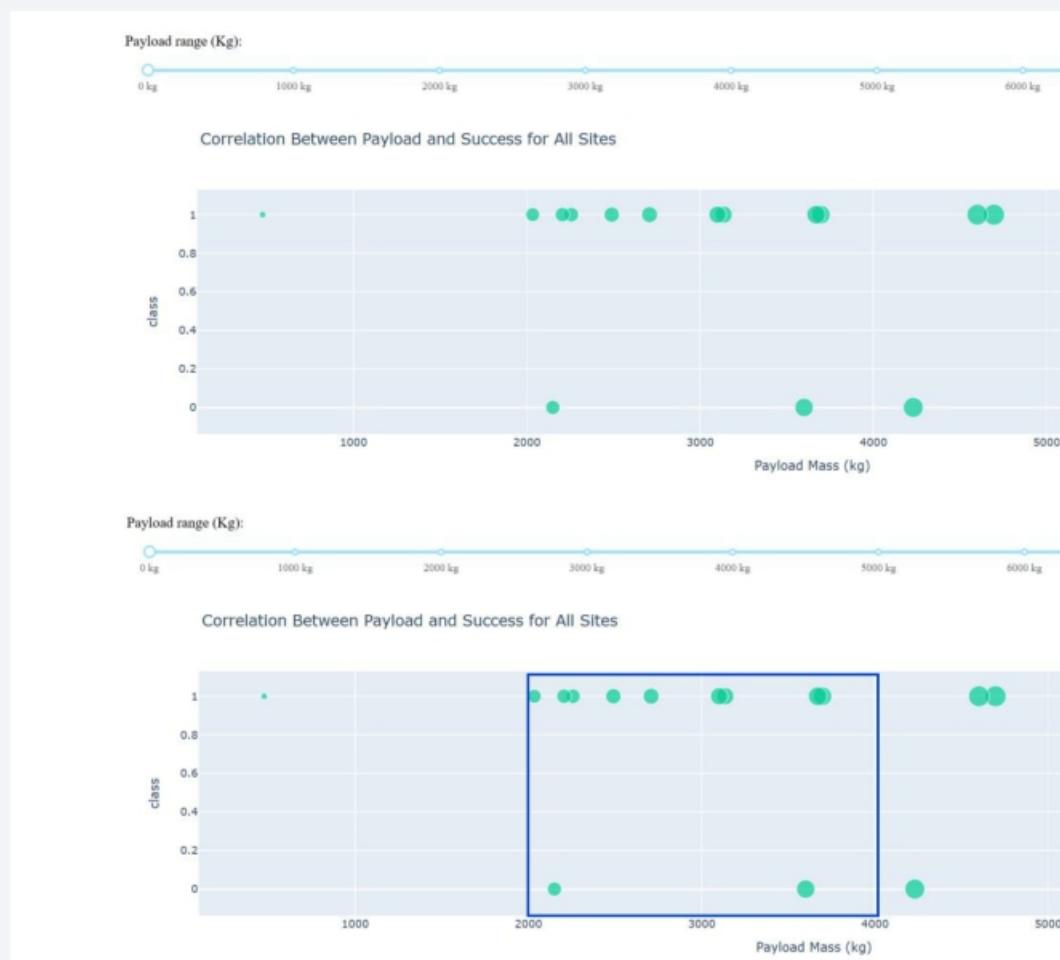
- **Launch Success Ratio:** Among all sites, the launch success was ranked third. 57.1% of all launches from CCAFS SLC-40 were successful.

Launch Outcomes By Payload



- **Launch Outcomes By Payload:** At a payload range of 2000 kg to 3200 kg, the success ratio is at 60.0%. In contrast, launch success ratio is 100% for payloads between 4000 kg and 6000 kg.

Launch Outcomes By Booster

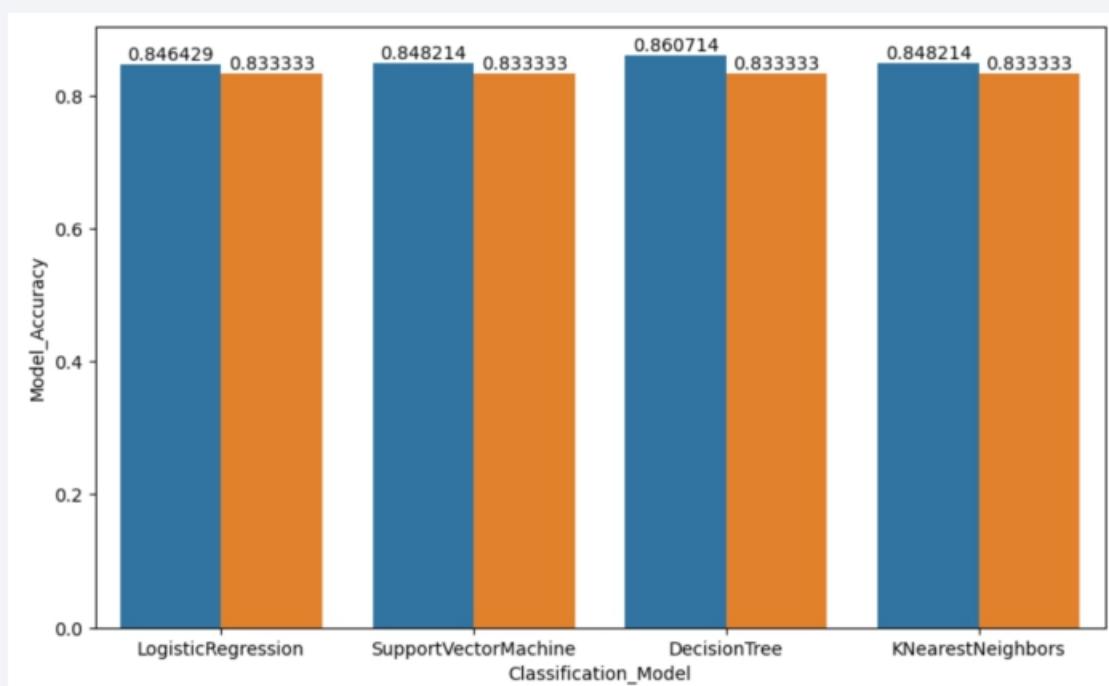


- **Launch Outcomes By Booster Version:** The highest launch success rate is achieved by the booster FT version at 65.0%. This increases to 81.8% when considering only successful launches.

Section 5

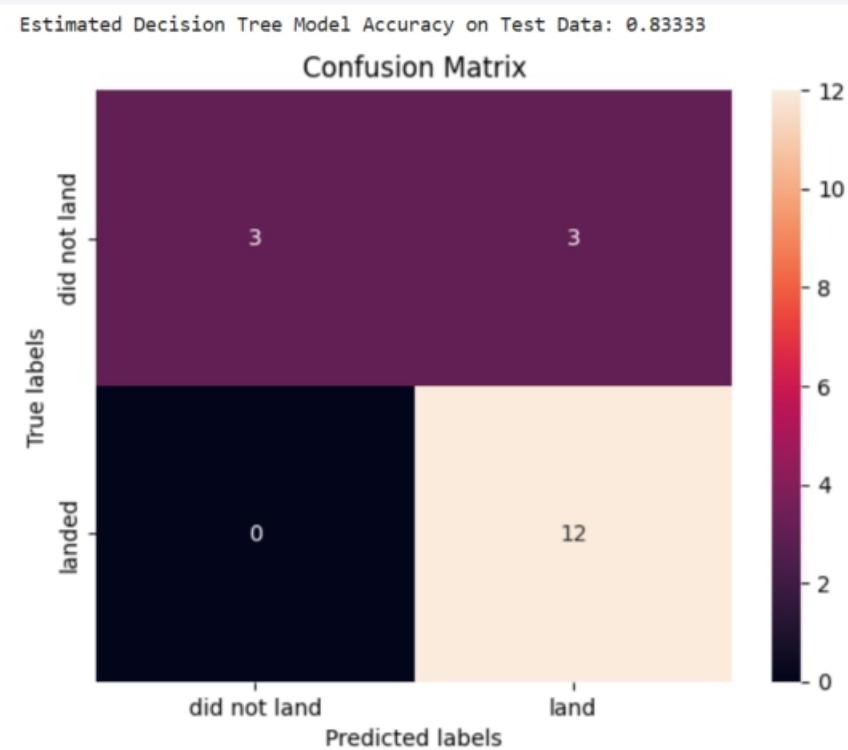
Predictive Analysis (Classification)

Classification Accuracy



- **Internal + External Model Evaluation:** All 4 candidate models achieved 83.33% accuracy on the test set. Higher accuracy was observed for the decision tree model at 86.07% based on internal evaluation using cross-validation.
- **Cross-Validated Accuracy on Train:** 86.07% (Decision Tree)
- **Accuracy on Test:** 83.33% (LR, SVM, Decision Tree, KNN)

Confusion Matrix



- **Classification Performance:** With a classification accuracy of 83.33% on the test set, the final selected model based on a decision tree correctly predicted 15 (12 true positive and 3 true negative cases) out of 18 total cases. 3 cases were noted which the model predicted as successful but did not actually land.

Conclusions

- Data collection for the analysis involved SpaceX REST API.
- Appropriate pre-processing methods including row and column deletion, missing value imputation, one-hot encoding, target creation were applied prior to subsequent analysis and modeling.
- EDA using visualization and SQL demonstrated the effect of independent variables such as flight number, launch site and orbit to successful landing outcomes.
- Interactive visual analytics using Folium demonstrated the effect of geographical proximities to a coastline, highways, railways and major cities to successful landing outcomes.
- Interactive visual analytics using Plotly Dash demonstrated the effect of independent variables such as launch site, payload mass and booster version to successful landing outcomes.
- A decision tree classification model provided a robust classification of landing outcomes with a cross-validated accuracy of 83.33%.

Thank you!