

Winning the Space Race with Data Science

James Porter
15 Aug 2023

Outline

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

Executive Summary

- **Summary of methodologies**
 - Data Collection
 - Data Wrangling
 - Exploratory Data Analysis with SQL
 - Exploratory Data Analysis Visualization
 - Visual Analytics with Folium
 - Interactive Dashboard with Plotly
 - Predictive Analysis
- **Summary of all results**
 - Exploratory Data Analysis Results
 - Interactive Analytics Demo
 - Predictive Analysis Results

Introduction

- Project Background and Context:
 - In this capstone, we will predict if the Falcon 9 first stage will land successfully. SpaceX advertises Falcon 9 rocket launches on its website, with a cost of 62 million dollars; other providers cost upward of 165 million dollars each, much of the savings is because SpaceX can reuse the first stage. Therefore if we can determine if the first stage will land, we can determine the cost of a launch. This information can be used if an alternate company wants to bid against SpaceX for a rocket launch.
- Questions wanting to answer:
 - Price per launch
 - Probability of a successful landing

Section 1

Methodology

Methodology

Executive Summary

- Data collection methodology:
 - SpaceX REST API
 - Web Scraping of Wikipedia
- Perform data wrangling
 - Filtering Data
 - Dealing with Missing Values
 - One Hot Encoding to Prepare Data
- Perform exploratory data analysis (EDA) using visualization and SQL

Methodology

Executive Summary

- Perform interactive visual analytics using Folium and Plotly Dash
- Perform predictive analysis using classification models
 - Collected data was normalized, divided into training and test data sets, and applied to four different classification models.

Data Collection

- THE PROCESS OF GATHERING DATA ENCOMPASSED A BLEND OF API REQUESTS FROM SPACEX's REST API AND THE EXTRACTION OF INFORMATION THROUGH WEB SCRAPING FROM A TABLE WITHIN SPACEX's WIKIPEDIA PAGE. THE UTILIZATION OF BOTH THESE METHODS WAS ESSENTIAL TO ACQUIRE A COMPREHENSIVE DATASET FOR A MORE INTRICATE ANALYSIS.
- THE DATA COLUMNS WERE POPULATED USING SPACEX's REST API, INCLUDING FLIGHtnumber, DATE, BOOSTERVERSION, PAYLOADMASS, ORBIT, LAUNCHSITE, OUTCOME, FLIGHTS, GRIDFINS, REUSED, LEGS, LANDINGPAD, BLOCK, REUSEDCOUNT, SERIAL, LONGITUDE, AND LATITUDE.
- FURTHERMORE, THE DATA COLUMNS SOURCED FROM WIKIPEDIA WEB SCRAPING ENCOMPASSED FLIGHT NUMBER, LAUNCH SITE, PAYLOAD, PAYLOADMASS, ORBIT, CUSTOMER, LAUNCH OUTCOME, BOOSTER VERSION, BOOSTER LANDING, DATE, AND TIME.

Data Collection – SpaceX API

- SPACEX PROVIDES AN API TO ALLOW ACCESS TO IT'S DATA.
- THIS DATA WAS COLLECTED AND PROCESSED TO USE AS PART OF OUR DATASET.
- [GITHUB – DATA COLLECTION API](#)

**Request the SpaceX API
and collect the launch
data**

**Filter Data to only show
Falcon 9 rocket data**

Deal with missing values

Data Collection - Scraping

- WIKIPEDIA ALSO CONTAINS A WEALTH OF INFORMATION ABOUT SPACEX LAUNCHES.
- GITHUB – DATA COLLECTION WEB SCRAPING

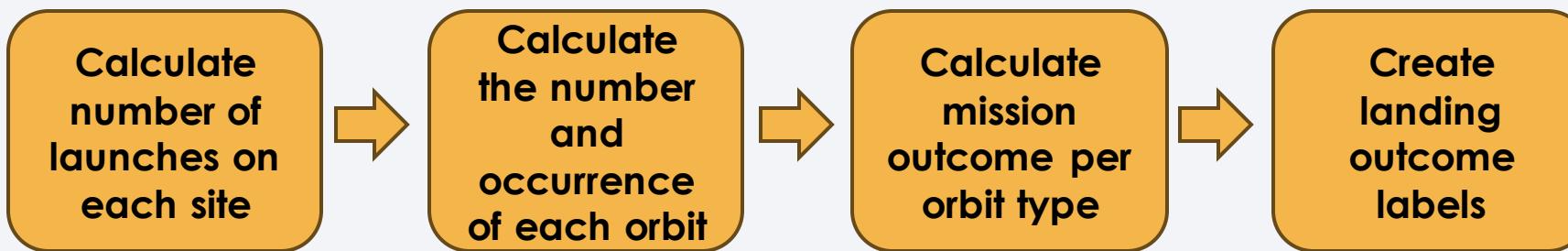
Locate the Falcon 9 launch data Wiki Page

Extract names from the HTML table header

Create a Data Frame from the collected HTML tables

Data Wrangling

- HERE THE DATA WAS EXPLORED TO IDENTIFY THE DIFFERENT TYPES OF DATA THAT WE COULD USE TO FORMULATE IDENTIFYING FACTORS THAT CONTRIBUTED TO EITHER SUCCESSFUL OR FAILED LANDING ATTEMPTS.



- [GITHUB – DATA WRANGLING](#)

EDA with Data Visualization

- THE FOLLOWING CHARTS WERE PLOTTED. THESE WERE CHOSEN TO SEE HOW DIFFERENT VARIABLES AFFECTED THE LANDING OUTCOME:
 - FLIGHTNUMBER VS. PAYLOADMASS
 - FLIGHTNUMBER VS LAUNCHSITE
 - PAYLOAD Vs. LAUNCHSITE
 - SUCESS RATE OF EACH ORBIT
 - FLIGHTNUMBER AND ORBIT TYPE.
 - PAYLOAD VS. ORBIT
- [GITHUB – EDA WITH DATA VISUALIZATION](#)

EDA with SQL

- SQL QUERIES PERFORMED:
 - `SELECT DISTINCT LAUNCH_SITE FROM SPACEXTBL`
 - `SELECT * FROM SPACEXTBL WHERE LAUNCH_SITE LIKE 'CCA%' LIMIT 5;`
 - `SELECT SUM(PAYLOAD_MASS_KG_) FROM SPACEXTBL WHERE CUSTOMER = 'NASA (CRS)'`
 - `SELECT AVG(PAYLOAD_MASS_KG_) FROM SPACEXTBL WHERE BOOSTER_VERSION LIKE 'F9 v1.1';`
 - `SELECT MIN(DATE) AS MIN_DATE FROM SPACEXTBL WHERE "LANDING_OUTCOME" = 'SUCCESS (GROUND PAD)';`
 - `SELECT BOOSTER_VERSION FROM SPACEXTBL WHERE (PAYLOAD_MASS_KG_ > 4000 AND PAYLOAD_MASS_KG_ < 6000)`
 - `AND (LANDING_OUTCOME = 'SUCCESS (DRONE SHIP)');`
 - `SELECT MISSION_OUTCOME, COUNT(MISSION_OUTCOME) AS COUNTS FROM SPACEXTBL GROUP BY MISSION_OUTCOME;`
 - `SELECT BOOSTER_VERSION, PAYLOAD_MASS_KG_ FROM SPACEXTBL WHERE PAYLOAD_MASS_KG_ = (SELECT MAX(PAYLOAD_MASS_KG_) FROM SPACEXTBL);`
 - `SELECT LANDING_OUTCOME, BOOSTER_VERSION, LAUNCH_SITE FROM SPACEXTBL WHERE LANDING_OUTCOME = 'FAILURE (DRONE SHIP)' AND YEAR(DATE) = '2015'`
 - `SELECT LANDING_OUTCOME, COUNT(*) AS LANDINGCOUNTS FROM SPACEXTBL WHERE DATE BETWEEN '2010-06-04' AND '2017-03-20'`
 - `GROUP BY LANDING_OUTCOME`
 - `ORDER BY COUNT(*) DESC;`
- [GITHUB – EDA WITH DATA VISUALIZATION](#)

Interactive Map with Folium

- MARKERS INDICATING LAUNCH SITES
 - ADDED BLUE CIRCLE AT NASA JOHNSON SPACE CENTER'S COORDINATE WITH A POPUP LABEL
 - SHOWING ITS NAME USING ITS LATITUDE AND LONGITUDE COORDINATES
 - ADDED RED CIRCLES AT ALL LAUNCH SITES COORDINATES WITH A POPUP LABEL
 - SHOWING ITS NAME USING ITS NAME USING ITS LATITUDE AND LONGITUDE COORDINATES
- COLORED MARKERS OF LAUNCH OUTCOMES
 - ADDED COLORED MARKERS OF SUCCESSFUL (GREEN) AND UNSUCCESSFUL (RED) LAUNCHES AT EACH LAUNCH SITE TO SHOW WHICH LAUNCH SITES HAVE HIGH SUCCESS RATES
- DISTANCES BETWEEN A LAUNCH SITE TO PROXIMITIES
 - ADDED COLORED LINES TO SHOW DISTANCE BETWEEN LAUNCH SITE CCAFS SLC-40 AND ITS PROXIMITY TO THE NEAREST COASTLINE, RAILWAY, HIGHWAY, AND CITY

Dashboard with Plotly Dash

- DROPODOWN LIST WITH LAUNCH SITES
 - ALLOW USER TO SELECT ALL LAUNCH SITES OR A CERTAIN LAUNCH SITE
- PIE CHART SHOWING SUCCESSFUL LAUNCHES
 - ALLOW USER TO SEE SUCCESSFUL AND UNSUCCESSFUL LAUNCHES AS A PERCENT OF THE TOTAL
- SLIDER OF PAYLOAD MASS RANGE
 - ALLOW USER TO SELECT PAYLOAD MASS RANGE
- SCATTER CHART SHOWING PAYLOAD MASS vs. SUCCESS RATE BY BOOSTER VERSION
 - ALLOW USER TO SEE THE CORRELATION BETWEEN PAYLOAD AND LAUNCH SUCCESS
- GITHUB – PLOTLY DASHBOARD APP

Predictive Analysis (Classification)

- WE LOADED THE DATA USING NUMPY AND PANDAS, TRANSFORMED THE DATA, SPLIT OUR DATA INTO TRAINING AND TESTING.
- WE BUILT DIFFERENT MACHINE LEARNING MODELS AND TUNE DIFFERENT HYPERPARAMETERS USING GRIDSEARCHCV.
- WE USED ACCURACY AS THE METRIC FOR OUR MODEL, IMPROVED THE MODEL USING FEATURE ENGINEERING AND ALGORITHM TUNING.
- WE FOUND THE BEST PERFORMING CLASSIFICATION MODEL.
- [GITHUB – PREDICTIVE ANALYSIS](#)

Results

Exploratory Data Analysis

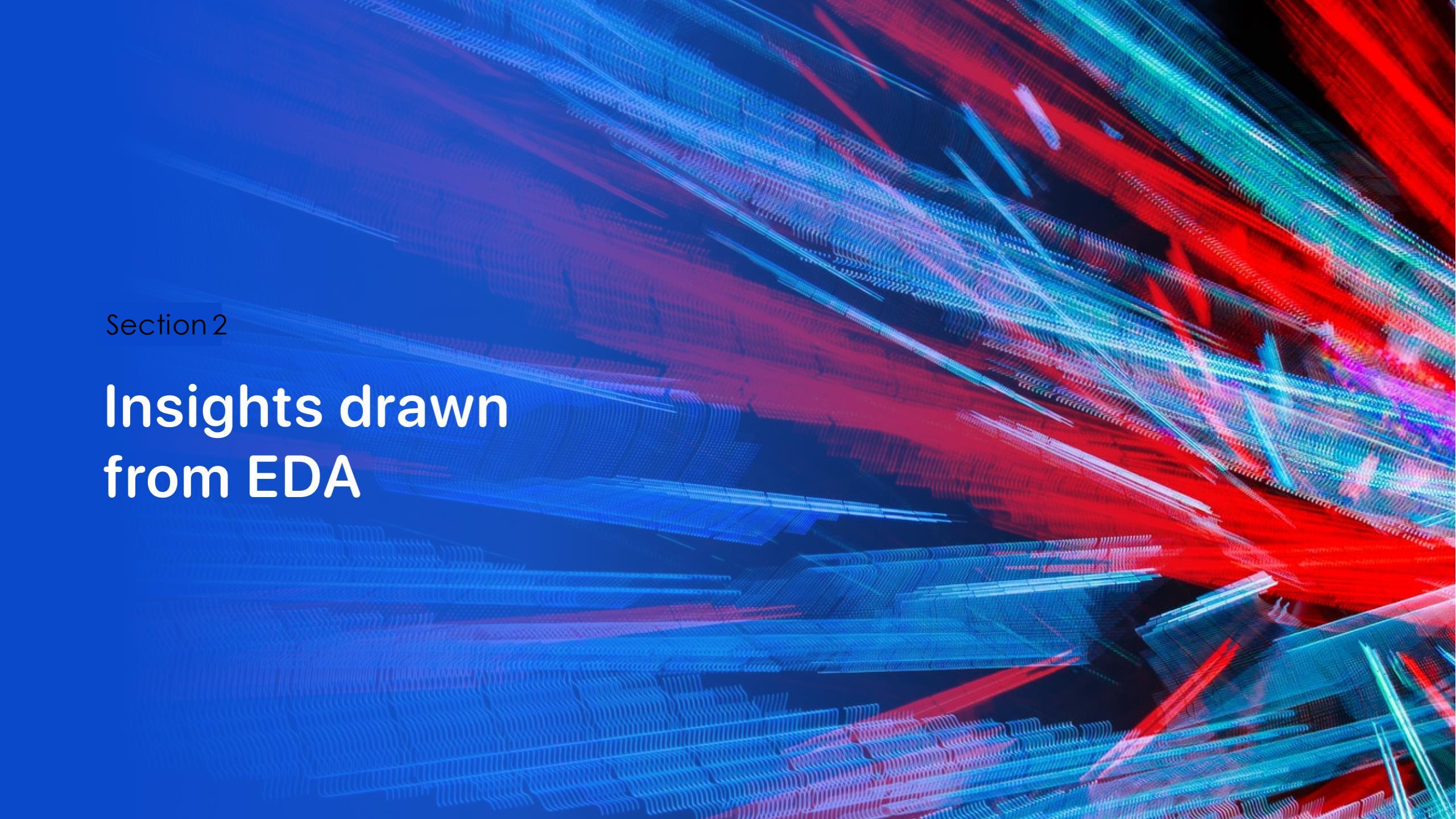
- Launch success has improved over time
- KSC LC-39A has the highest success rate among landing sites
- Orbits ES-L1, GEO, HEO and SSO have a 100% success rate

Visual Analytics

- Most launch sites are near the equator, and all are close to the coast
- Launch sites are far enough away from anything a failed launch can damage
- (city, highway, railway), while still close enough to bring people and material
- to support launch activities

Predictive Analytics

- Decision Tree model is the best predictive model for the dataset

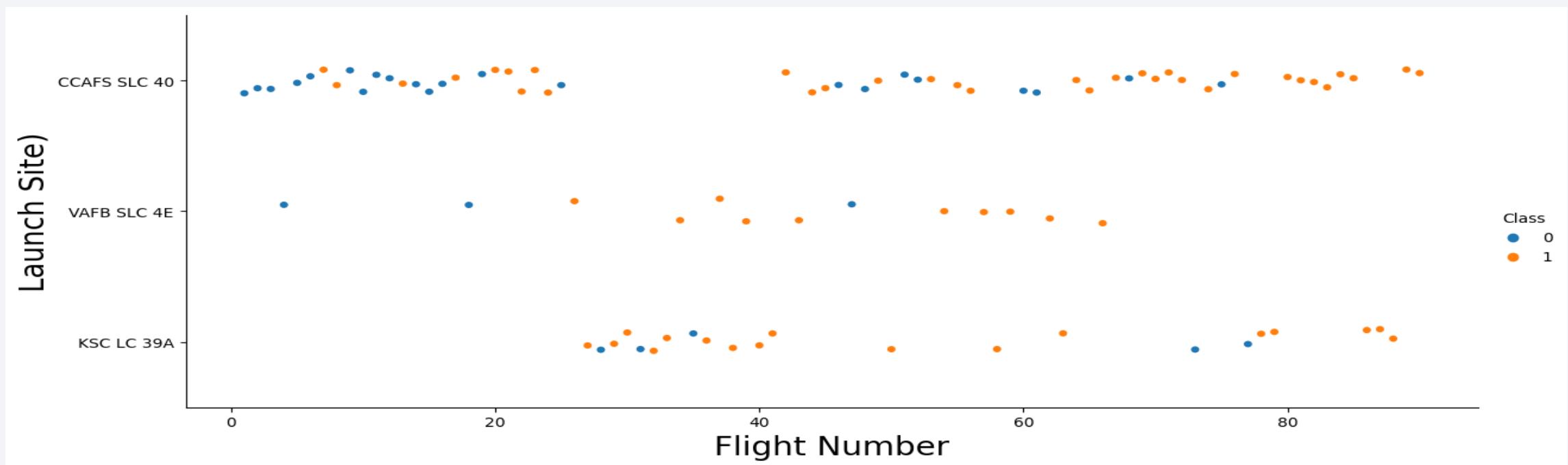
The background of the slide features a complex, abstract digital visualization. It consists of numerous thin, glowing lines that create a sense of depth and motion. The lines are primarily blue and red, with some green and purple highlights. They form a grid-like structure that curves and twists across the frame, resembling a 3D wireframe or a network of data points. The overall effect is futuristic and dynamic, suggesting concepts like data flow, digital communication, or complex systems.

Section 2

Insights drawn from EDA

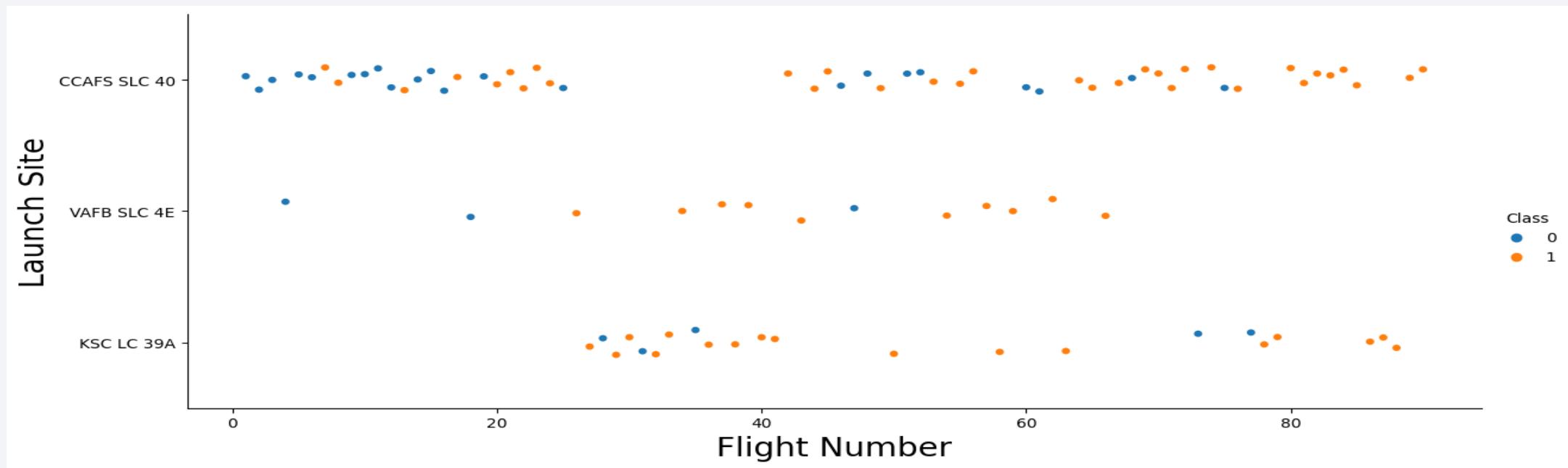
Flight Number vs. Launch Site

- FROM THE PLOT, WE FOUND THAT THE MORE FLIGHTS AT A LAUNCH SITE, THE
- GREATER THE SUCCESS RATE AT A LAUNCH SITE.



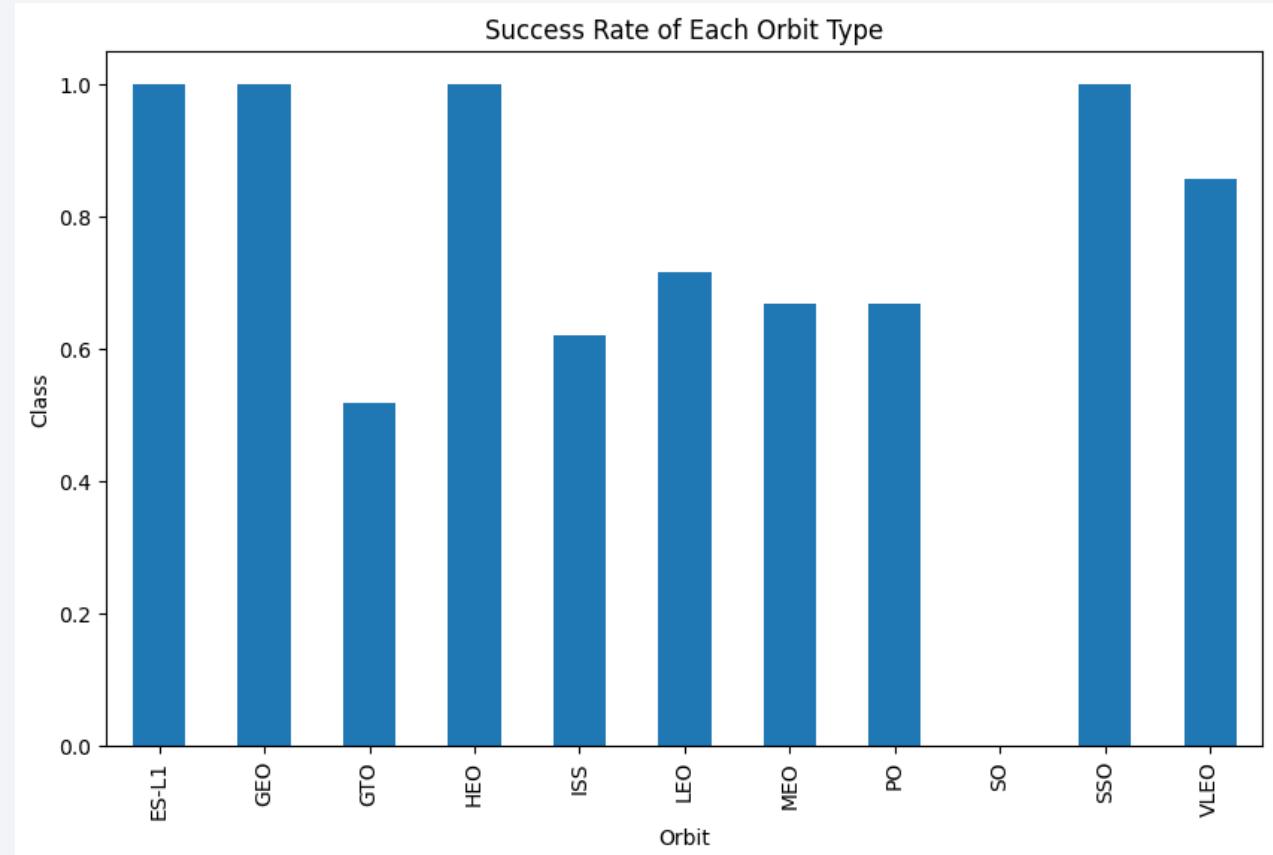
Payload vs. Launch Site

- THE GREATER THE PAYLOAD, THE MORE SUCCESSFUL THE LAUNCH WAS AT SITE CCAFS SLC-40



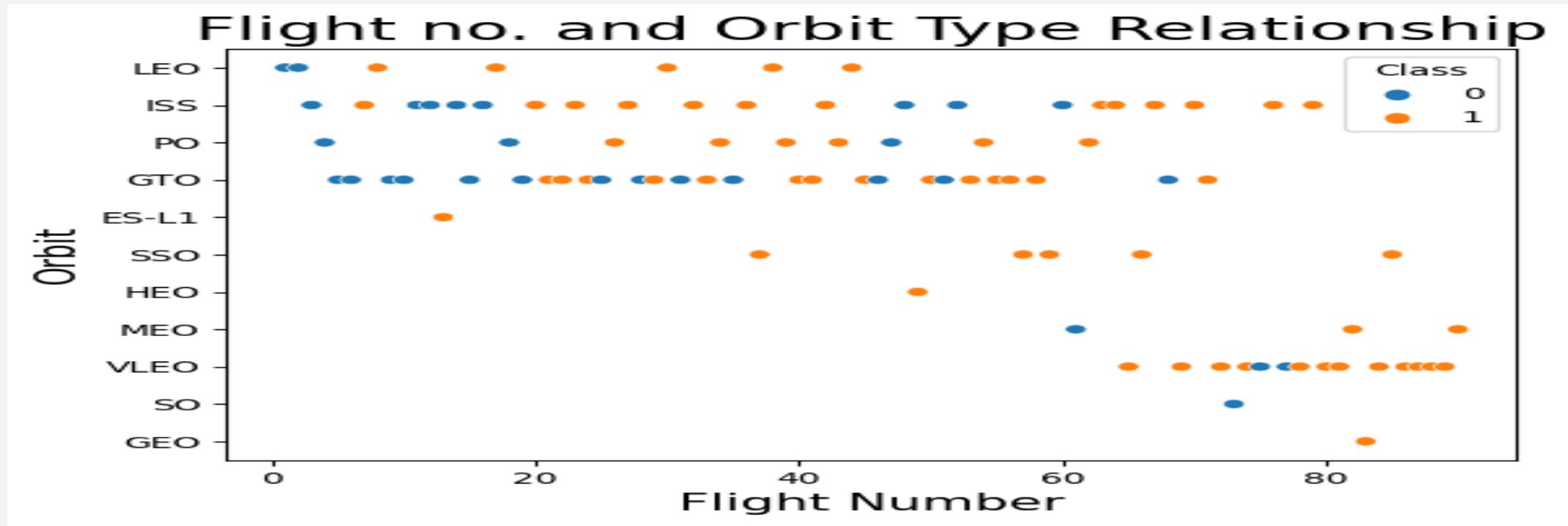
Success Rate vs. Orbit Type

- ES-L1, GEO, HEO,
- SSO, VLEO HAD THE HIGHEST
- SUCCESS RATE.



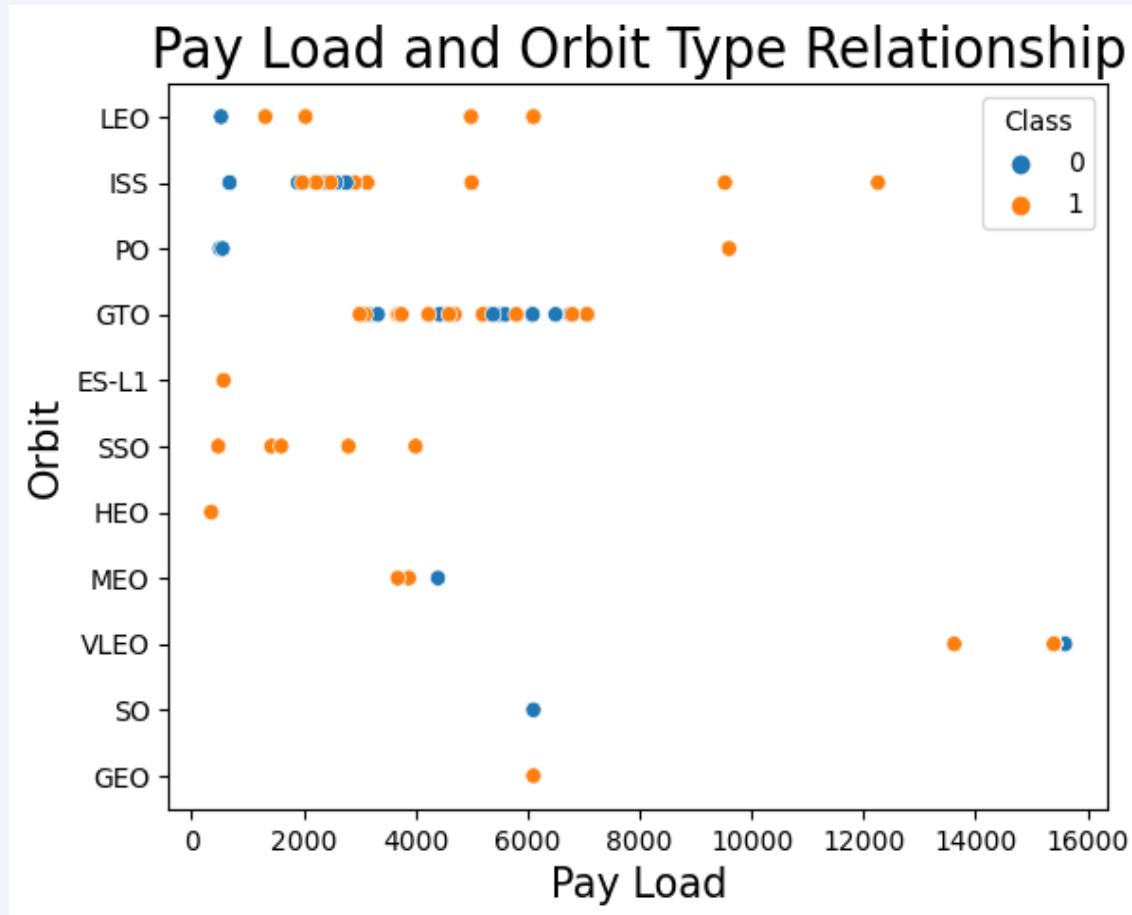
Flight Number vs. Orbit Type

- WE SEE THAT IN THE LEO ORBIT, SUCCESS IS RELATED TO THE NUMBER OF FLIGHTS WHEREAS IN THE GTO ORBIT, THERE IS NO RELATIONSHIP BETWEEN FLIGHT NUMBER AND THE ORBIT.



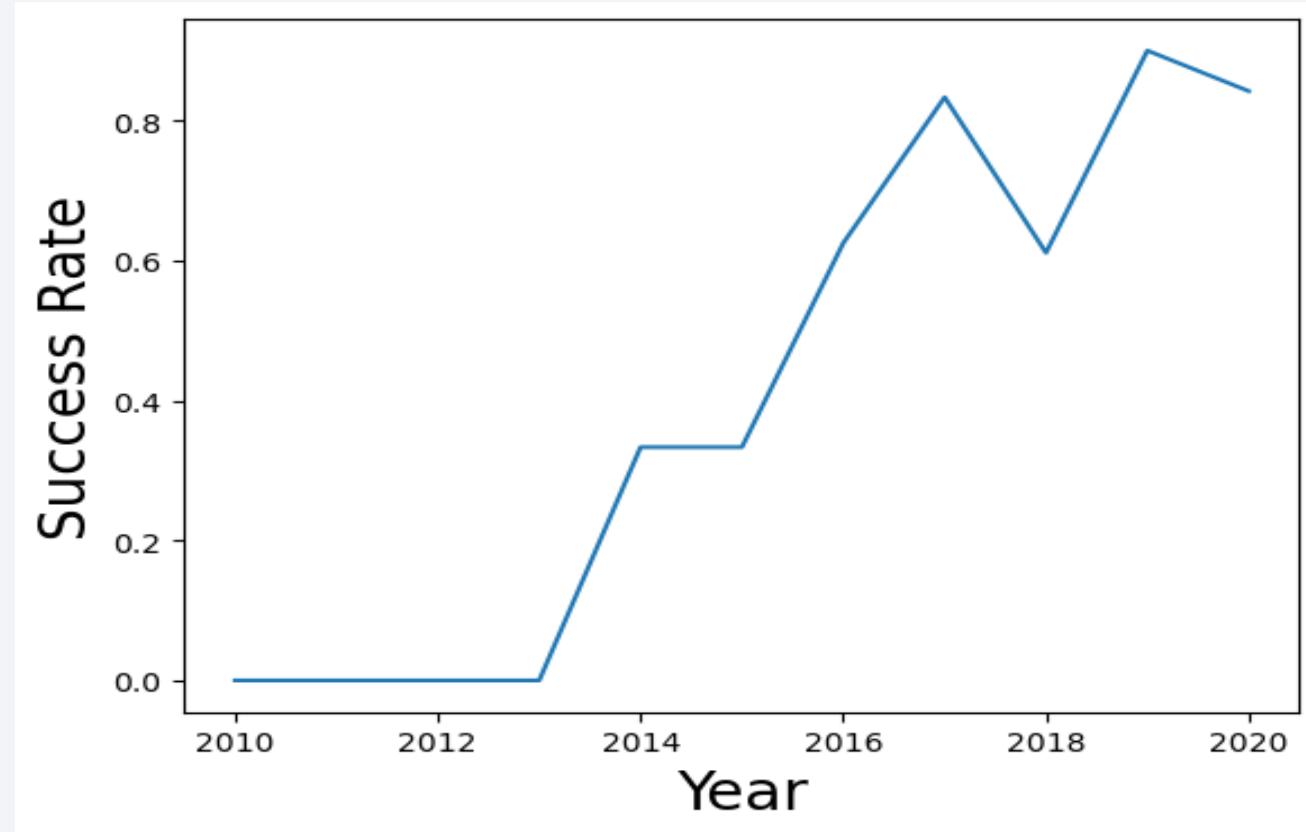
Payload vs. Orbit Type

- PO, LEO AND ISS ORBITS HAD MORE SUCCESSFUL LANDINGS WITH A HEAVIER PAYLOAD.



Launch Success Yearly Trend

- WE CAN SEE THE LAUNCH SUCCESS HAS INCREASED OVER TIME.



All Launch Site Names

- WE USED THE KEY WORD **DISTINCT** TO SHOW ONLY
- UNIQUE LAUNCH SITES FROM THE SPACEX DATA.

```
%%sql
select distinct Launch_Site from spacextbl
* sqlite:///my_data1.db
Done.



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


```

Launch Site Names Begin with 'CCA'

- WE USED WHERE TO DISPLAY 5 RECORDS WHERE LAUNCH SITES BEGIN WITH `CCA`

SpaceX Launch Record Selection										
Date	Time (UTC)	Booster_Version	Launch_Site	Payload	PAYLOAD_MASS_KG_	Orbit	Customer	Mission_Outcome	Landing_Outcome	
2010-04-06	18:45:00	F9 v1.0 B0003	CCAFS LC-40	Dragon Spacecraft Qualification Unit	0	LEO	SpaceX	Success	Failure (parachute)	
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	Failure (parachute)	
2012-05-22	07:44:00	F9 v1.0 B0005	CCAFS LC-40	Dragon demo flight C2	525	LEO (ISS)	NASA (COTS)	Success	No attempt	
2012-08-10	00:35:00	F9 v1.0 B0006	CCAFS LC-40	SpaceX CRS-1	500	LEO (ISS)	NASA (CRS)	Success	No attempt	
2013-01-03	15:10:00	F9 v1.0 B0007	CCAFS LC-40	SpaceX CRS-2	677	LEO (ISS)	NASA (CRS)	Success	No attempt	

Total Payload Mass

- WE CALCULATED THE TOTAL PAYLOAD USING SUM AND WHERE.

```
%%sql

select sum(PAYLOAD_MASS__KG_) from spacextbl where Customer = 'NASA (CRS)'

* sqlite:///my_data1.db
Done.

sum(PAYLOAD_MASS__KG_)

45596
```

Average Payload Mass by F9 v1.1

- WE CALCULATED THE AVERAGE
- PAYLOAD MASS WITH AVG AND WHERE

```
%%sql
select avg(PAYLOAD_MASS_KG_) from spacextbl where Booster_Version LIKE 'F9 v1.1';
* sqlite:///my_data1.db
Done.
avg(PAYLOAD_MASS_KG_)
_____
2928.4
```

First Successful Ground Landing Date

- DATES OF THE FIRST SUCCESSFUL LANDING DATE MIN AND WHERE

```
%%sql
select min(Date) as min_date from spacextbl 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

- WE FOUND THAT THE DATES OF THE FIRST SUCCESSFUL LANDING OUTCOME ON
- GROUND PAD WAS 22ND DECEMBER 2015

```
%%sql
select Booster_Version from spacextbl where (PAYLOAD_MASS_KG_ > 4000 and PAYLOAD_MASS_KG_ < 6000)
and (Landing_Outcome = 'Success (drone ship)');
* 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

- THE TOTAL NUMBER OF SUCCESSFUL AND FAILURE MISSION OUTCOMES WERE FOUND COUNT AND FROM

```
%%sql
select Mission_Outcome, count(Mission_Outcome) as counts from spacextbl group by Mission_Outcome;
* sqlite:///my_data1.db
Done.



| Mission_Outcome                  | counts |
|----------------------------------|--------|
| Failure (in flight)              | 1      |
| Success                          | 98     |
| Success                          | 1      |
| Success (payload status unclear) | 1      |


```

Boosters Carried Maximum Payload

- THE NAMES OF THE BOOSTER WHICH HAVE CARRIED THE MAXIMUM PAYLOAD MASS WERE FOUND USING MAX

```
%%sql
select Booster_Version, PAYLOAD_MASS__KG_ from spacextbl where PAYLOAD_MASS__KG_ = (select max(PAYLOAD_MASS__KG_) from spacextbl)
* sqlite:///my_data1.db
Done.

Booster_Version    PAYLOAD_MASS__KG_
F9 B5 B1048.4      15600
F9 B5 B1049.4      15600
F9 B5 B1051.3      15600
F9 B5 B1056.4      15600
F9 B5 B1048.5      15600
F9 B5 B1051.4      15600
F9 B5 B1049.5      15600
F9 B5 B1060.2      15600
F9 B5 B1058.3      15600
F9 B5 B1051.6      15600
F9 B5 B1060.3      15600
F9 B5 B1049.7      15600
```

2015 Launch Records

- THE FAILED LANDING_OUTCOMES IN DRONE SHIP, THEIR BOOSTER VERSIONS, AND LAUNCH SITE NAMES FOR IN YEAR 2015 USING LIKE AND BETWEEN.

List the failed landing_outcomes in drone ship, their booster versions,

```
task_9 = """
    SELECT BoosterVersion, LaunchSite, LandingOutcome
    FROM SpaceX
    WHERE LandingOutcome LIKE 'Failure (drone ship)'
        AND Date BETWEEN '2015-01-01' AND '2015-12-31'
    ...
create_pandas_df(task_9, database=conn)
```

	boosterversion	launchsite	landingoutcome
0	F9 v1.1 B1012	CCAFS LC-40	Failure (drone ship)
1	F9 v1.1 B1015	CCAFS LC-40	Failure (drone ship)

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

- THE COUNT OF LANDING OUTCOMES (SUCH AS FAILURE (DRONE SHIP) OR SUCCESS (GROUND PAD)) BETWEEN THE DATE 2010-06-04 AND 2017-03-20, RANKED IN DESCENDING ORDER USING COUNT, GROUP BY, AND ORDER BY

```
%%sql
select Landing_Outcome, count(*) as LandingCounts from spacextbl where Date between '2010-06-04' and '2017-03-20'
group by Landing_Outcome
order by count(*) desc;

* sqlite:///my_data1.db
Done.



| Landing_Outcome        | LandingCounts |
|------------------------|---------------|
| 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             |


```

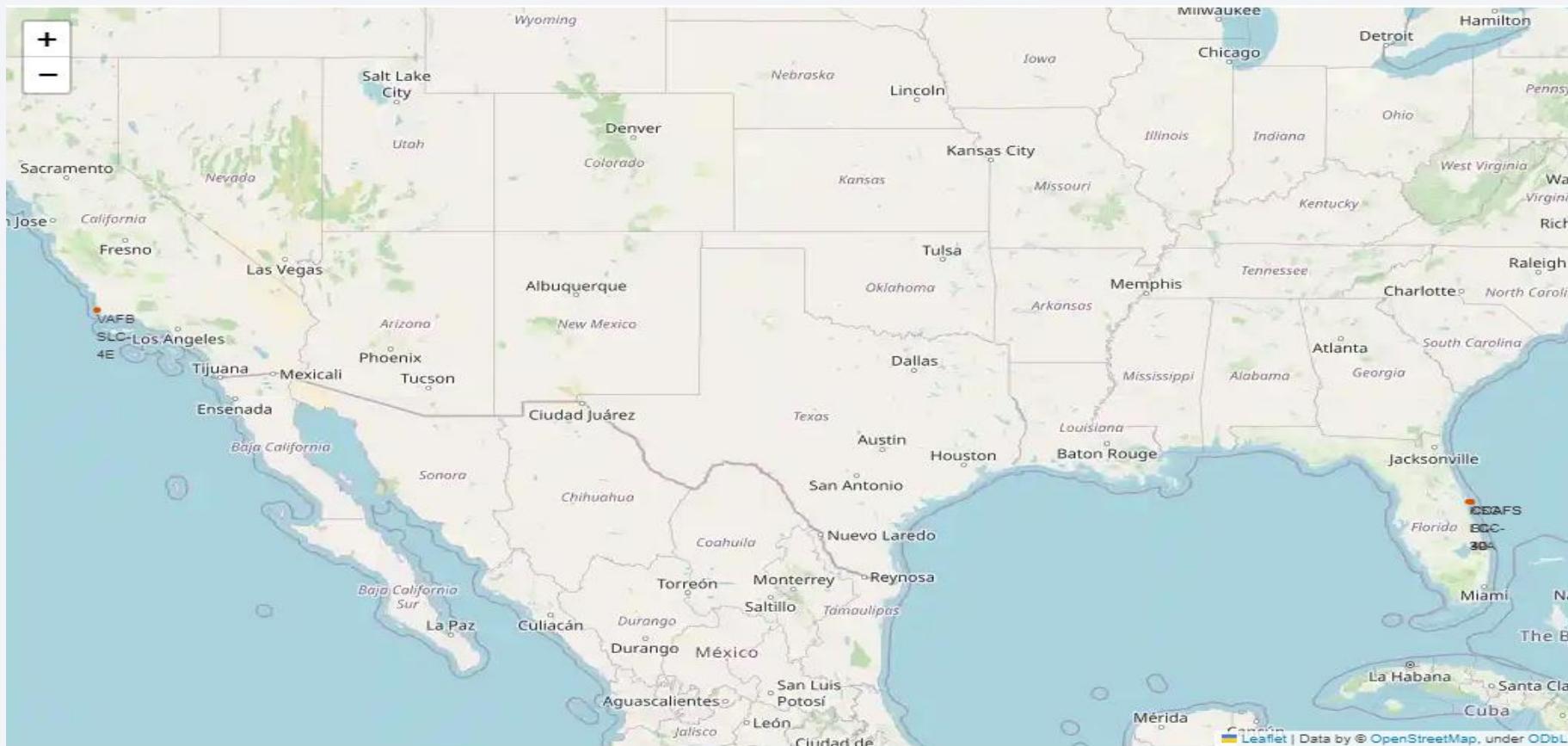
The background of the slide is a photograph taken from space at night. It shows the curvature of the Earth's horizon against a dark blue sky. City lights are visible as numerous small white and yellow dots, primarily concentrated in the lower right quadrant where the United States appears. In the upper right, there are bright green and yellow bands of light, likely the Aurora Borealis or Australis. The overall atmosphere is dark and mysterious.

Section 3

Launch Sites Proximities Analysis

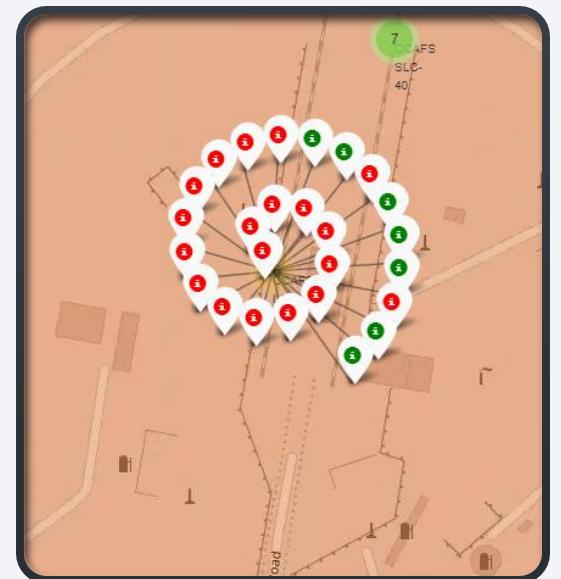
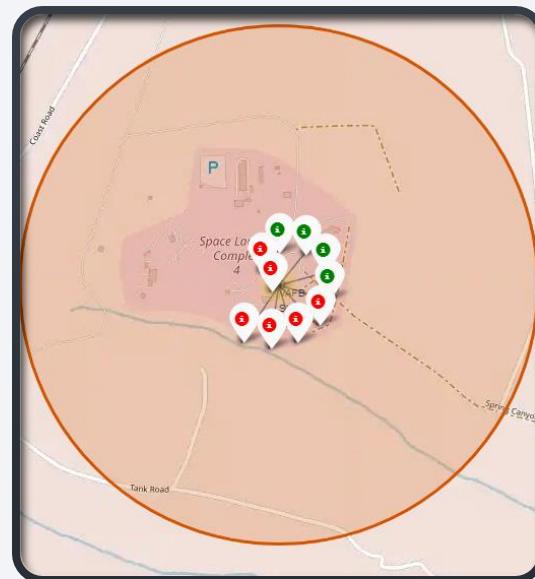
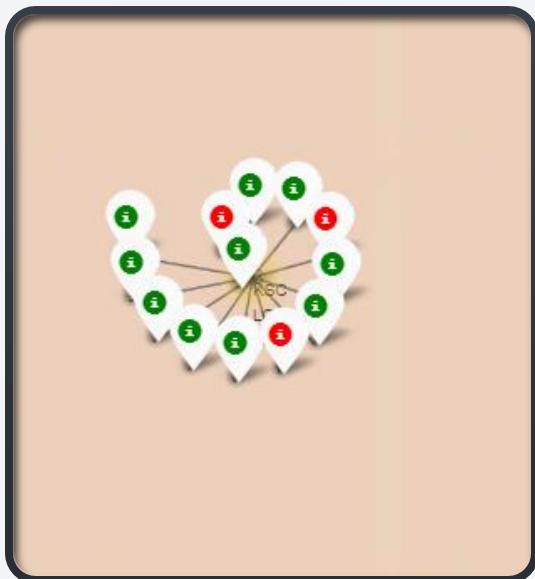
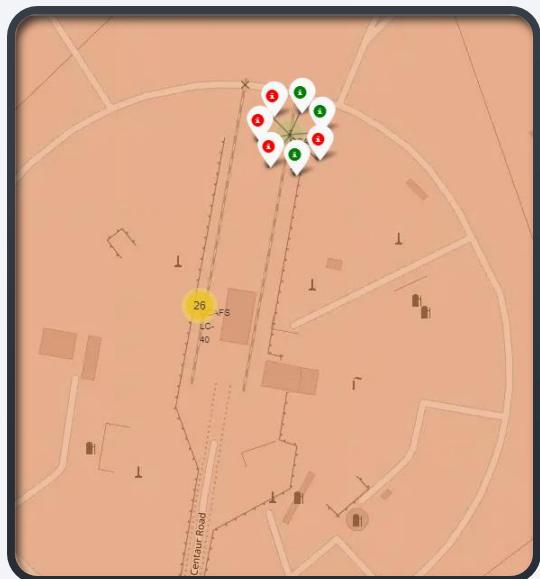
Launch Site Locations

- LAUNCH LOCATIONS ARE ON THE EAST COAST IN FLORIDA, AND THE WEST COAST IN CALIFORNIA.



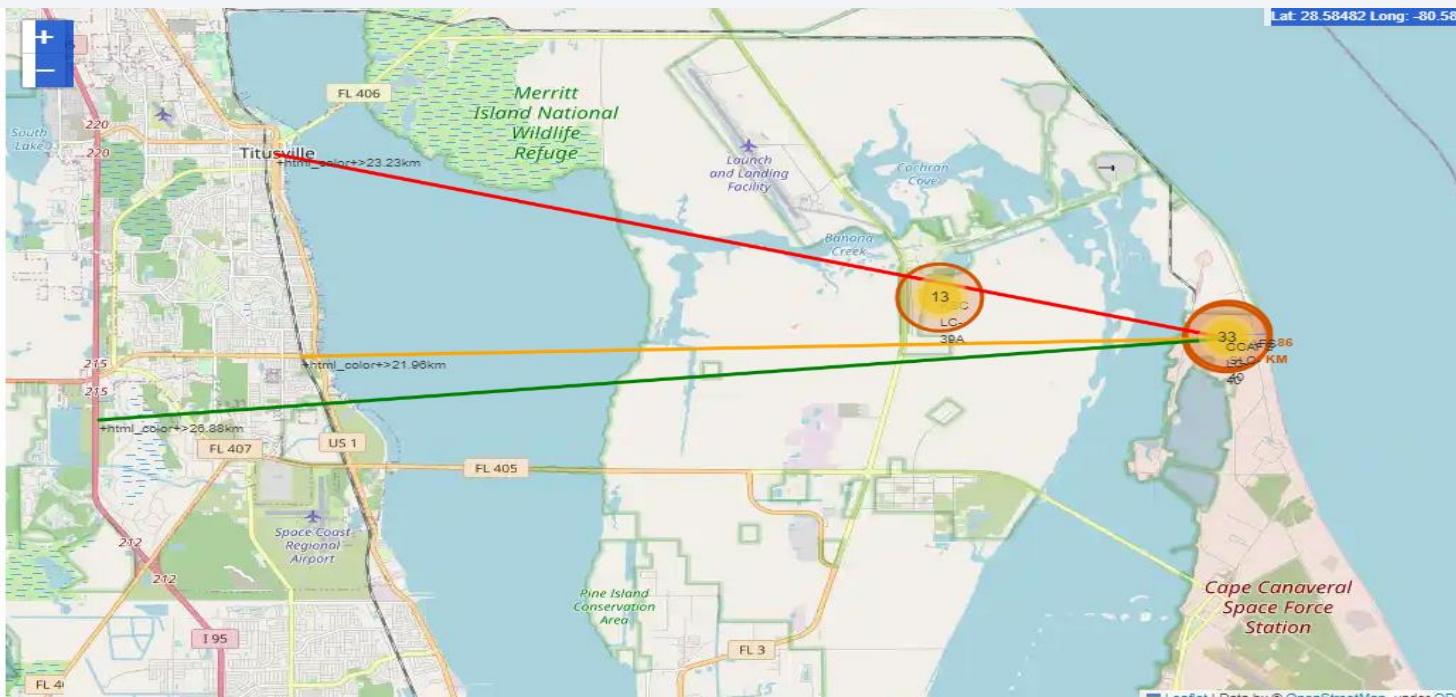
Markers Showing Launch Sites

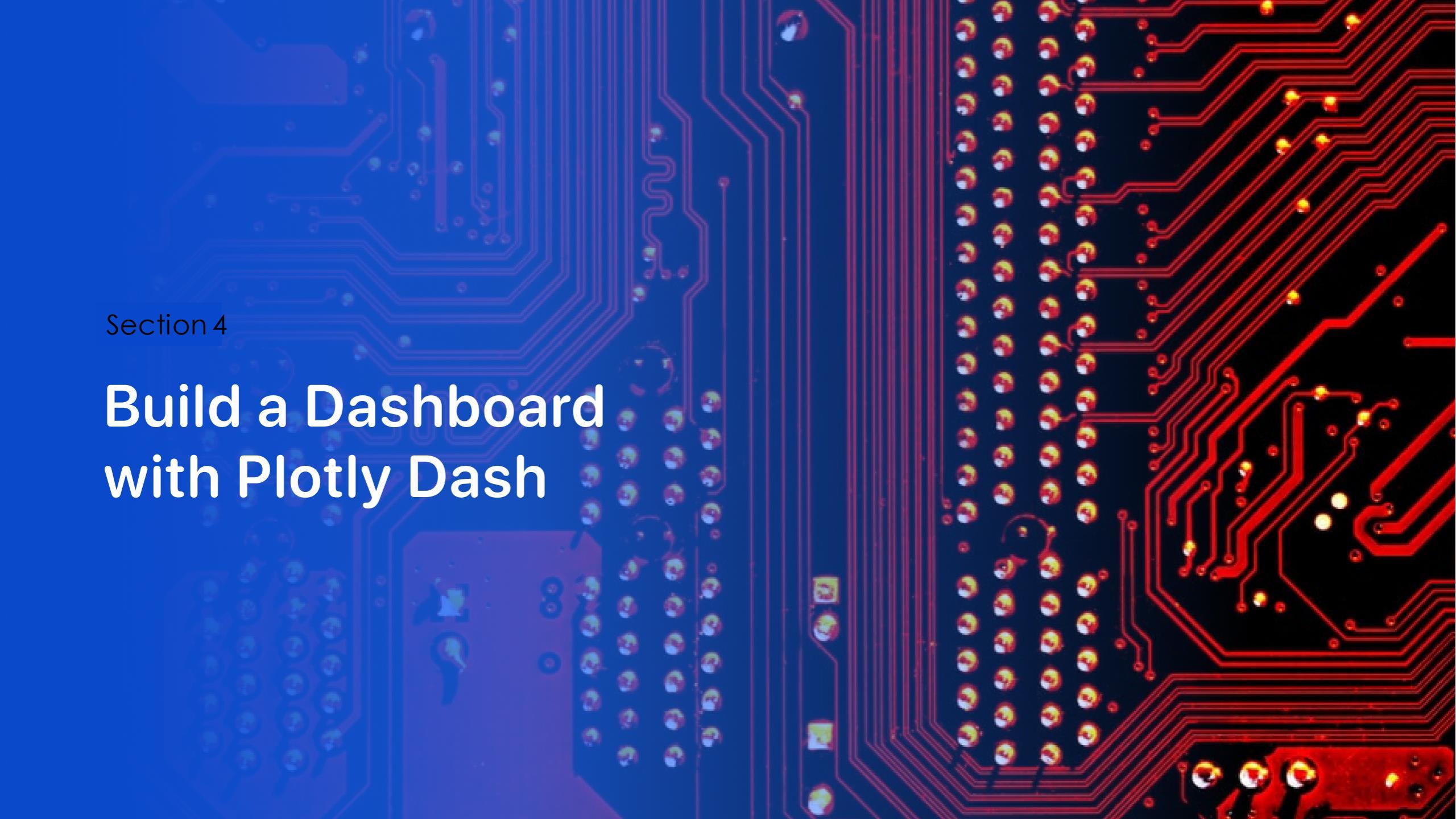
- SITES WITH COLOR CODED MARKERS INDICATING SUCCESSFUL/FAILED LAUNCHES PER SITE.



Pertinent Proximities

- MAP SHOWING PROXIMITIES TO CLOSEST RAILROAD, HIGHWAYS, AND CITY. INDICATES THAT CARE HAS BEEN TAKEN TO PICK LAUNCH SITES WITH A SAFE DISTANCE FROM OCCUPIED LOCATIONS



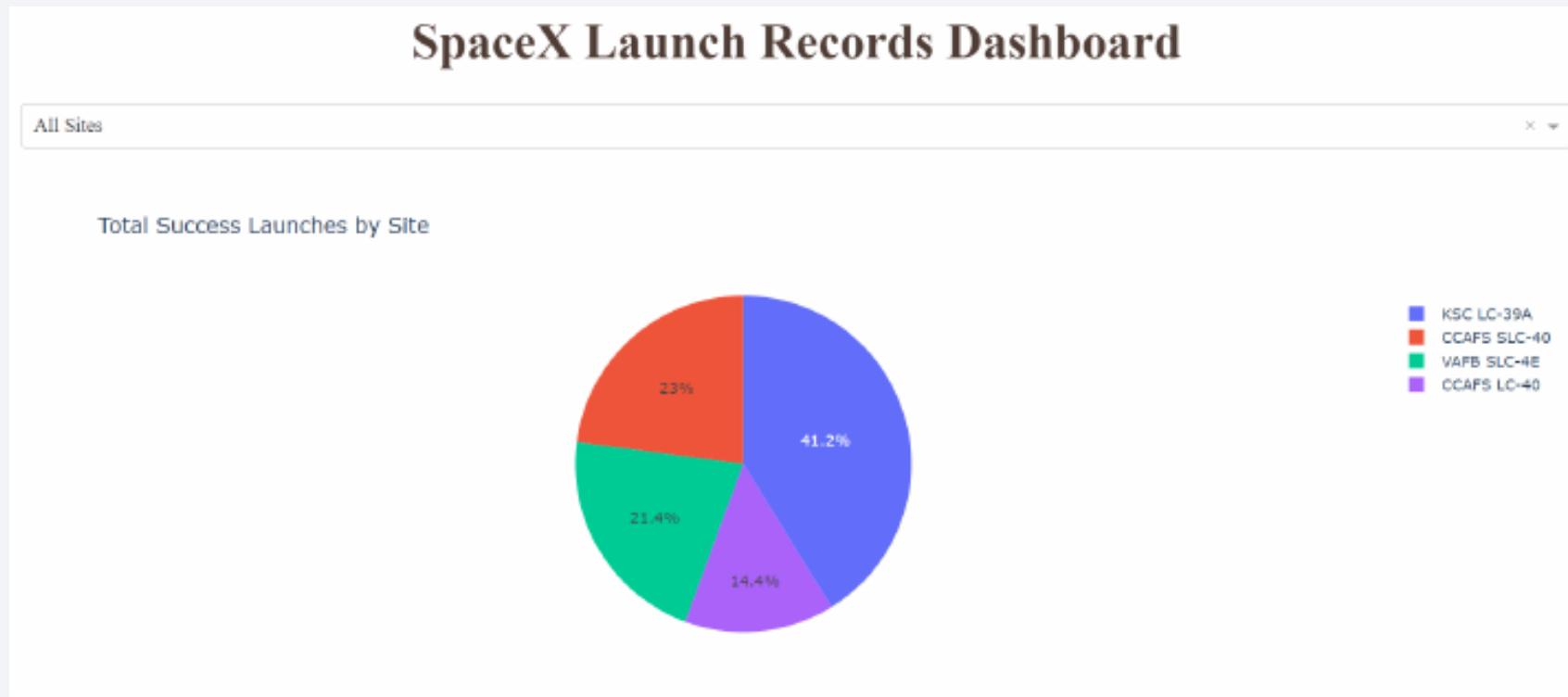


Section 4

Build a Dashboard with Plotly Dash

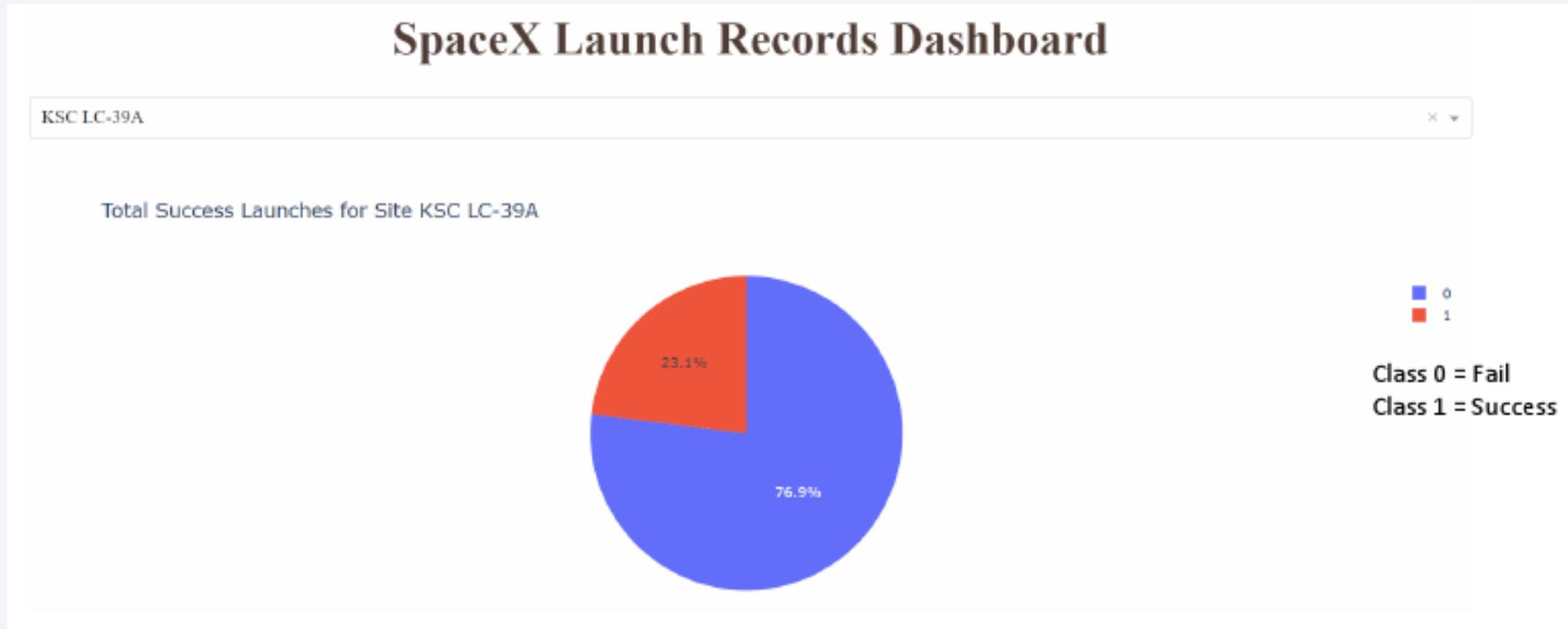
Successful Launches by Site

- KSC LC-39A HAS THE MOST SUCCESSFUL LAUNCHES AMONGST LAUNCH SITES (41.2%)



KSC LC-39A Successful Launches

- KSC LC-39A HAD THE HIGHEST SUCCESS RATE AMONG ALL LAUNCH SITES



Payload Vs Success for all Sites

- SCATTER PLOT OF PAYLOAD VS LAUNCH OUTCOME FOR ALL SITES



The background of the slide features a dynamic, abstract design. It consists of several thick, curved lines in shades of blue and yellow, creating a sense of motion and depth. The lines curve from the bottom left towards the top right, with some lines being more prominent than others. The overall effect is reminiscent of a tunnel or a high-speed journey through a digital space.

Section 5

Predictive Analysis (Classification)

Classification Accuracy

- THE DECISION TREE CLASSIFIER IS THE MODEL WITH THE HIGHEST CLASSIFICATION ACCURACY

```
# Create a DF for algorithm type and respective best scores  
  
Model_Performance_df = pd.DataFrame({'Algo Type': ['Logistic Regression', 'SVM','Decision Tree','KNN'],  
'Accuracy Score': [logreg_cv.best_score_, svm_cv.best_score_, tree_cv.best_score_, knn_cv.best_score_],  
'Test Data Accuracy Score': [logreg_cv.score(X_test, Y_test), svm_cv.score(X_test, Y_test),  
tree_cv.score(X_test, Y_test), knn_cv.score(X_test, Y_test)]})
```

```
Model_Performance_df.sort_values(['Accuracy Score'], ascending = False, inplace=True)
```

```
Model_Performance_df
```

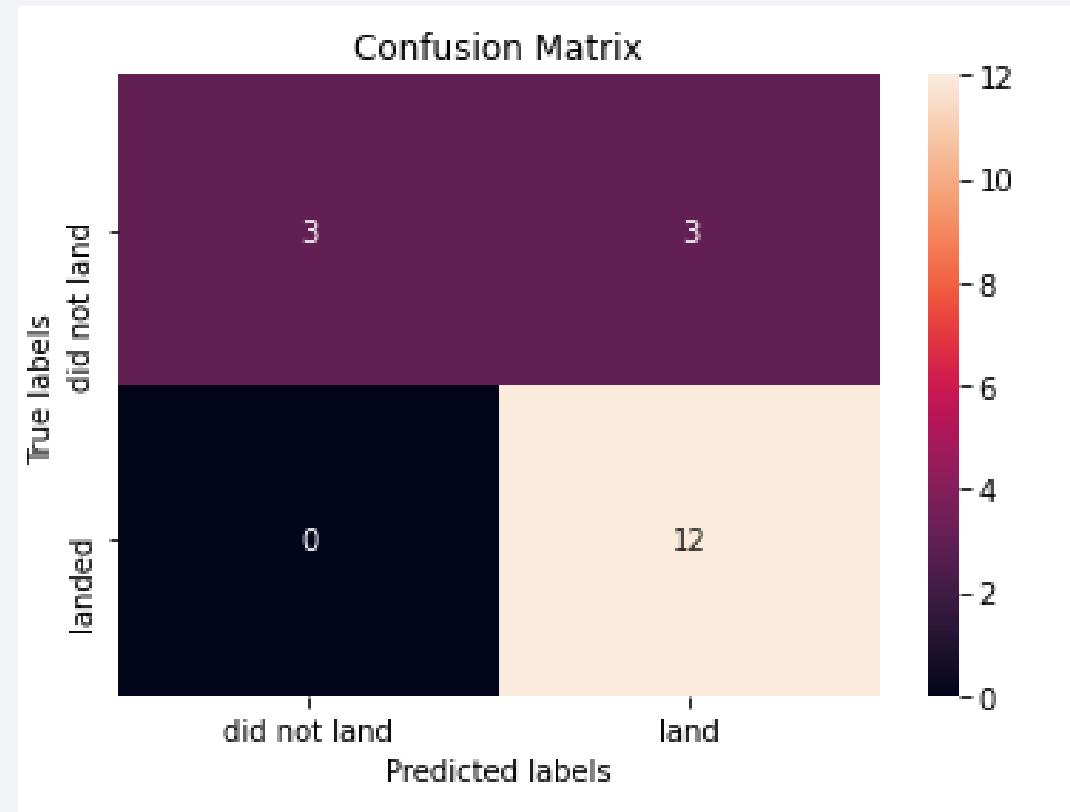
	Algo Type	Accuracy Score	Test Data Accuracy Score
2	Decision Tree	0.875000	0.833333
3	KNN	0.848214	0.833333
1	SVM	0.848214	0.833333
0	Logistic Regression	0.846429	0.833333

```
i = Model_Performance_df['Accuracy Score'].idxmax()  
print('The best performing alogrithm is '+ Model_Performance_df['Algo Type'][i]  
+ ' with score ' + str(Model_Performance_df['Accuracy Score'][i]))
```

```
The best performing alogrithm is Decision Tree with score 0.875
```

Confusion Matrix

- THE CONFUSION MATRIXES FOR ALL THE MODELS WERE THE SAME. FALSE POSITIVES WERE A PROBLEM IN ALL MODELS.



Conclusions

- WE CAN CONCLUDE THAT:
- THE LARGER THE FLIGHT AMOUNT AT A LAUNCH SITE, THE GREATER THE SUCCESS RATE AT A LAUNCH SITE.
- LAUNCH SUCCESS RATE STARTED TO INCREASE IN 2013 TILL 2020.
- ORBITS ES-L1, GEO, HEO, SSO, VLEO HAD THE MOST SUCCESS RATE.
- KSC LC-39A HAD THE MOST SUCCESSFUL LAUNCHES OF ANY SITES.
- THE DECISION TREE CLASSIFIER IS THE BEST MACHINE LEARNING ALGORITHM FOR THIS TASK.

A photograph of a night sky over a beach. A bright, white rocket launch is visible, curving upwards from the horizon. The sky is filled with stars and a mix of dark and illuminated clouds. In the foreground, a person walks along the wet sand near the water's edge. The text "Thank You!" is overlaid in the upper right quadrant.

Thank You!