

Winning Space Race with Data Science

ROHAN MALICK

10/8/2024

<https://github.com/coderRohan123>



Outline



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

Executive Summary

- Collected data from public SpaceX API and SpaceX Wikipedia page. Created labels column ‘class’ which classifies successful landings. Explored data using SQL, visualization, folium maps, and dashboards. Gathered relevant columns to be used as features. Changed all categorical variables to binary using one hot encoding. Standardized data and used GridSearchCV to find best parameters for machine learning models. Visualize accuracy score of all models.
- Four machine learning models were produced: Logistic Regression, Support Vector Machine, Decision Tree Classifier, and K Nearest Neighbors. All produced similar results with accuracy rate of about 83.33%. All models over predicted successful landings. More data is needed for better model determination and accuracy.

Introduction



SpaceX Falcon 9 Rocket – The Verge

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

Section 1

Methodology

Methodology

- **Data collection methodology:**
 - Combined data from SpaceX public API and SpaceX Wikipedia page
- **Perform data wrangling**
 - Classifying true landings as successful and unsuccessful otherwise
- **Perform exploratory data analysis (EDA) using visualization and SQL**
- **Perform interactive visual analytics using Folium and Plotly Dash**
- **Perform predictive analysis using classification models**
 - Tuned models using GridSearchCV

Data Collection Overview

Data collection process involved a combination of API requests from Space X public API and web scraping data from a table in Space X's Wikipedia entry.

The next slide will show the flowchart of data collection from API and the one after will show the flowchart of data collection from webscraping.

Space X API Data Columns:

FlightNumber, Date, BoosterVersion, PayloadMass, Orbit, LaunchSite, Outcome, Flights, GridFins,
Reused, Legs, LandingPad, Block, ReusedCount, Serial, Longitude, Latitude

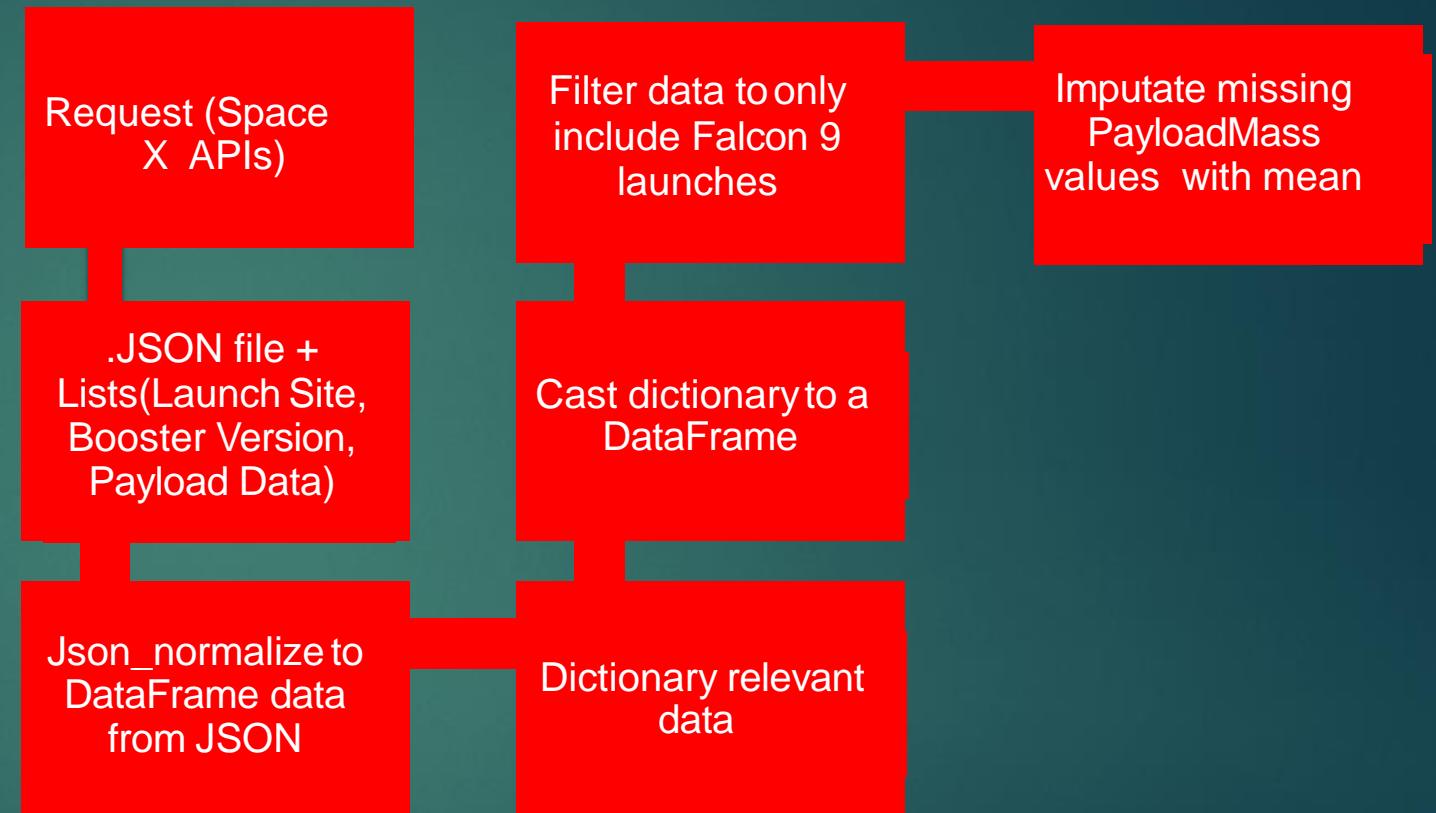
Wikipedia Webscrape Data Columns:

Flight No., Launch site, Payload, PayloadMass, Orbit, Customer, Launch outcome, Version Booster, Booster landing, Date, Time

Data Collection— SpaceX API

GitHub url:

<https://github.com/coderRohan123/IBM-DATASCIENCE-PROFESSIONAL-CERTIFICATE/blob/main/capstone/1.1data-collection-api.ipynb>



Data Collection – Web Scraping

GitHub url:

<https://github.com/coderRohan123/IBM-DATASCIENCE-PROFESSIONAL-CERTIFICATE/blob/main/capstone/1.2lab-s-webscraping.ipynb>

Request
Wikipedia
html

BeautifulSou
p
html5lib
Parser

Find launch
info html
table

Cast dictionary
to
DataFrame

Iterate
through table
cells to
extract data
to dictionary

Create
dictionary

Data Wrangling

- ▶ Create a training label with landing outcomes where successful = 1 & failure = 0.
- ▶ Outcome column has two components: ‘Mission Outcome’ ‘Landing Location’
- ▶ New training label column ‘class’ with a value of 1 if ‘Mission Outcome’ is True and 0 otherwise.
- ▶ Value Mapping:
- ▶ True ASDS, True RTLS, & True Ocean – set to -> 1
- ▶ None None, False ASDS, None ASDS, False Ocean, False RTLS – set to -> 0
- ▶ GitHub url:

<https://github.com/coderRohan123/IBM-DATASCIENCE-PROFESSIONAL-CERTIFICATE/blob/main/capstone/1.3-Data%20wrangling.ipynb>

EDA with Data Visualization

Exploratory Data Analysis performed on variables Flight Number, Payload Mass, Launch Site, Orbit, Class and Year.

Plots Used:

Flight Number vs. Payload Mass, Flight Number vs. Launch Site, Payload Mass vs. Launch Site, Orbit vs. Success Rate, Flight Number vs. Orbit, Payload vs Orbit, and Success Yearly Trend

Scatter plots, line charts, and bar plots were used to compare relationships between variables to decide if a relationship exists so that they could be used in training the machine learning model

GitHub url:

<https://github.com/coderRohan123/IBM-DATASCIENCE-PROFESSIONAL-CERTIFICATE/blob/main/capstone/2.1%20Exploratory%20Data%20Analysis%20-%20SQL.ipynb>

Loaded data set into IBM DB2 Database.

Queried using SQL Python integration.

Queries were made to get a better understanding of the dataset.

Queried information about launch site names, mission outcomes, various pay load sizes of customers and booster versions, and landing outcomes

GitHub url:

<https://github.com/coderRohan123/IBM-DATASCIENCE-PROFESSIONAL-CERTIFICATE/blob/main/capstone/2.2%20eda%20data%20viz.ipynb>

Build an interactive map with Folium

13

Folium maps mark Launch Sites, successful and unsuccessful landings, and a proximity example to key locations: Railway, Highway, Coast, and City.

This allows us to understand why launch sites may be located where they are. Also visualizes successful landings relative to location.

GitHub url:

https://github.com/coderRohan123/IBM-DATASCIENCE-PROFESSIONAL-CERTIFICATE/blob/main/capstone/3.1%20launch_site_location.ipynb

Build a Dashboard with Plotly Dash

Dashboard includes a pie chart and a scatter plot.

Pie chart can be selected to show distribution of successful landings across all launch sites and can be selected to show individual launch site success rates.

Scatter plot takes two inputs: All sites or individual site and payload mass on a slider between 0 and 10000 kg.

The pie chart is used to visualize launch site success rate.

The scatter plot can help us see how success varies across launch sites, payload mass, and booster version category.

GitHub url:

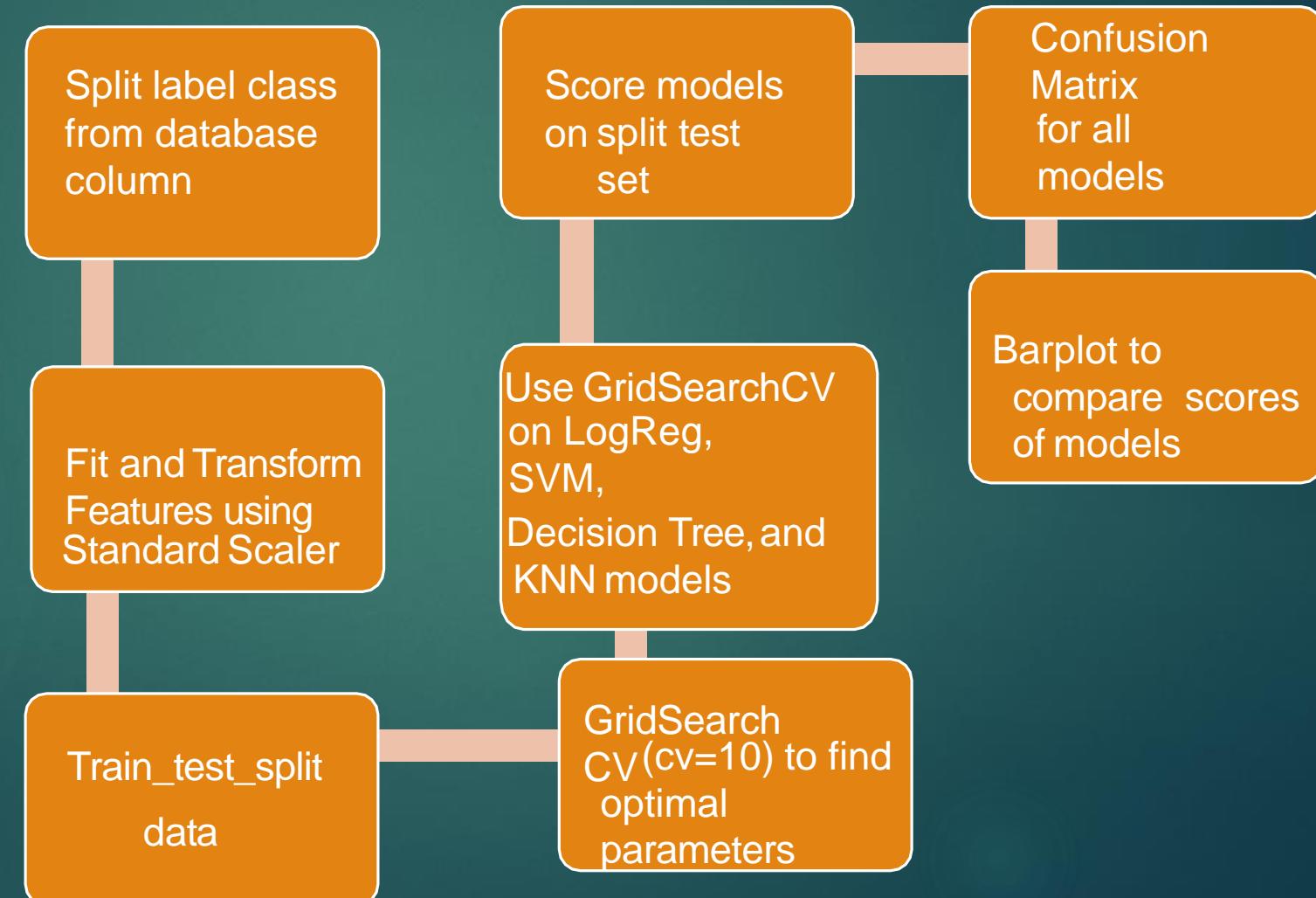
https://github.com/coderRohan123/IBM-DATASCIENCE-PROFESSIONAL-CERTIFICATE/blob/main/capstone/spacex_dash_app.py

Predictive analysis (Classification)

15

GitHub url:

https://github.com/coderRohan123/IBM-DATASCIENCE-PROFESSIONAL-CERTIFICATE/blob/main/capstone/4.1%20SpaceX_Machine%20Learning%20Prediction_Part_5.ipynb



Results

Exploratory Data Analysis

Interactive Analytics

Predictive Analysis

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 three-dimensional space or a network of data points. The overall effect is futuristic and dynamic.

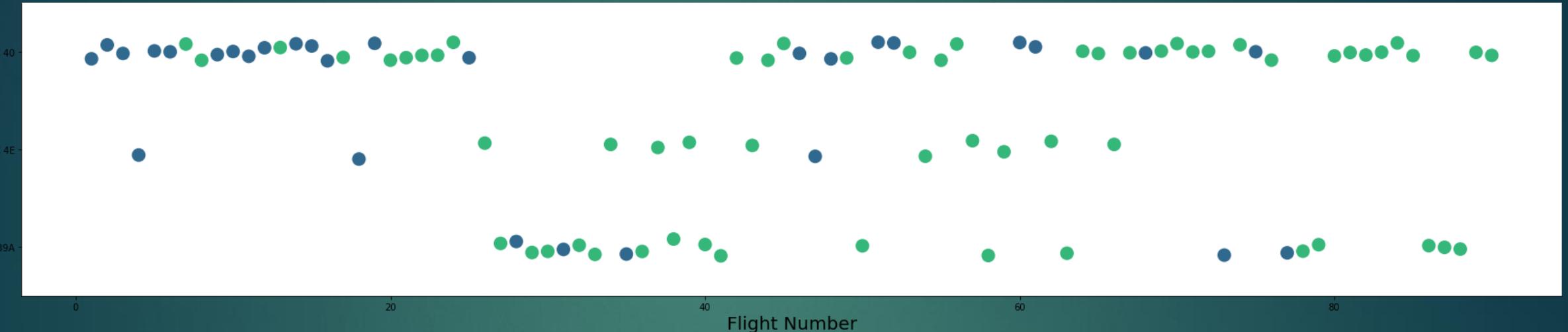
Section 2

Insights drawn from EDA

EDA - WITH VISUALIZATION

Flight Number vs. LaunchSite

19

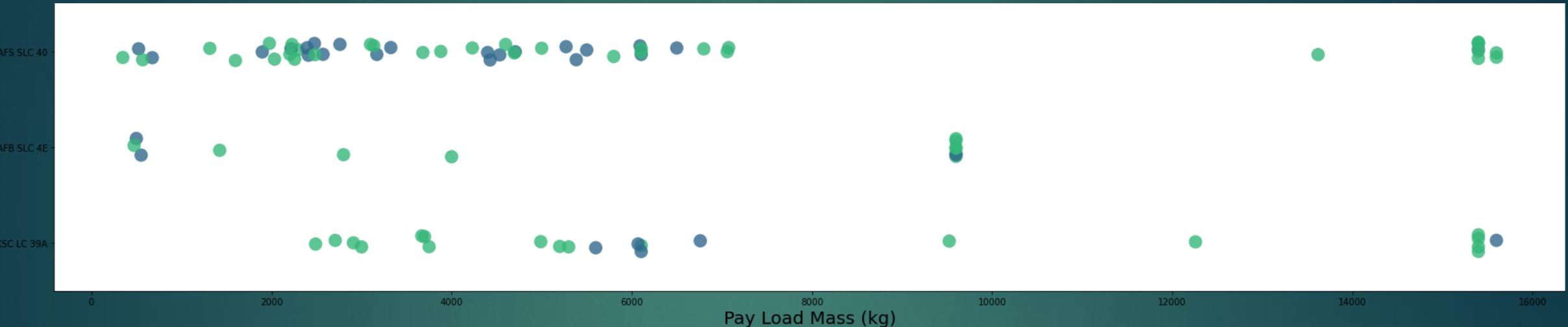


Green indicates successful launch; Purple indicates unsuccessful launch.

- Graphic suggests an increase in success rate over time (indicated in Flight Number).
- Likely a big breakthrough around flight 20 which significantly increased success rate.
- CCAFS appears to be the main launch site as it has the most volume.

Payload vs. Launch Site

20

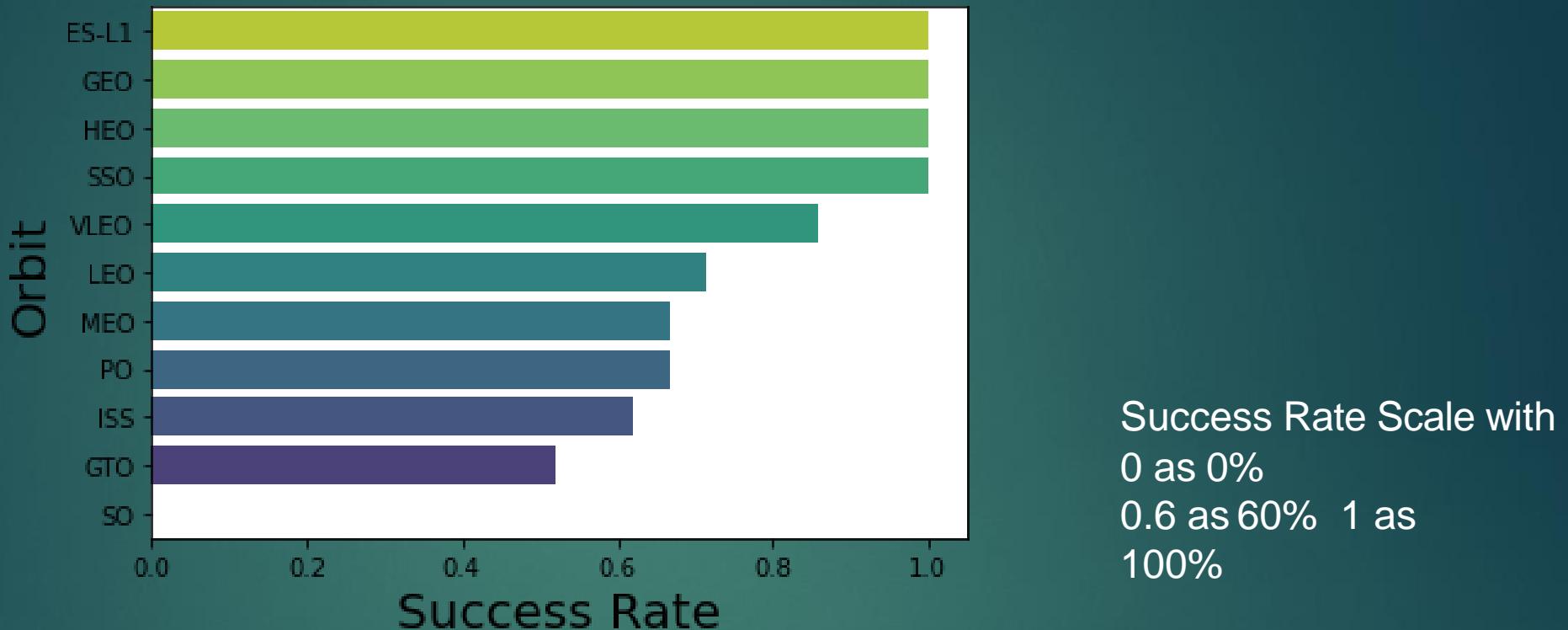


Green indicates successful launch; Purple indicates unsuccessful launch.

- Payload mass appears to fall mostly between 0-6000 kg
- Different launch sites also seem to use different payload mass

Success rate vs. Orbittype

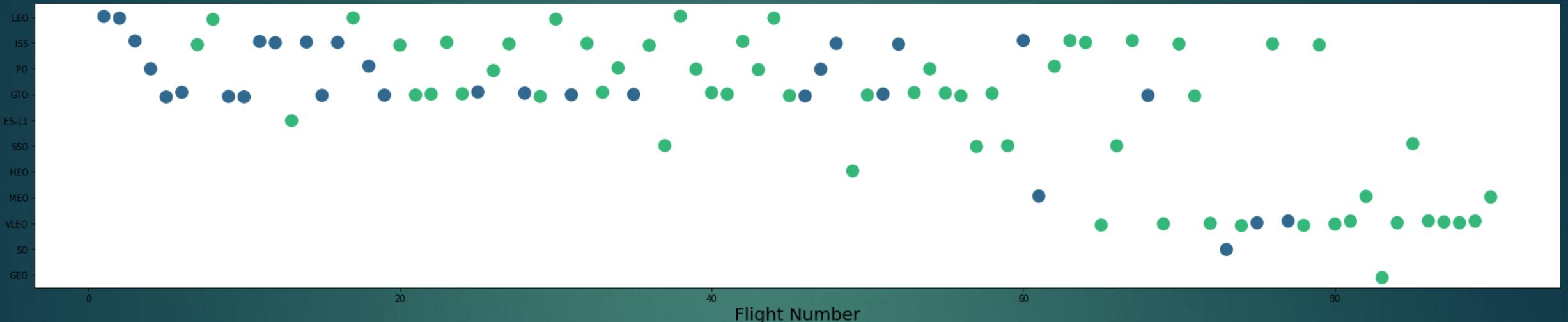
21



- ES-L1 (1), GEO (1), HEO (1) have 100% success rate (sample sizes in parenthesis) SSO (5) has 100% success rate
- VLEO (14) has decent success rate and attempts
- SO (1) has 0% success rate
- GTO (27) has the around 50% success rate but largest sample

Flight Number vs. Orbittype

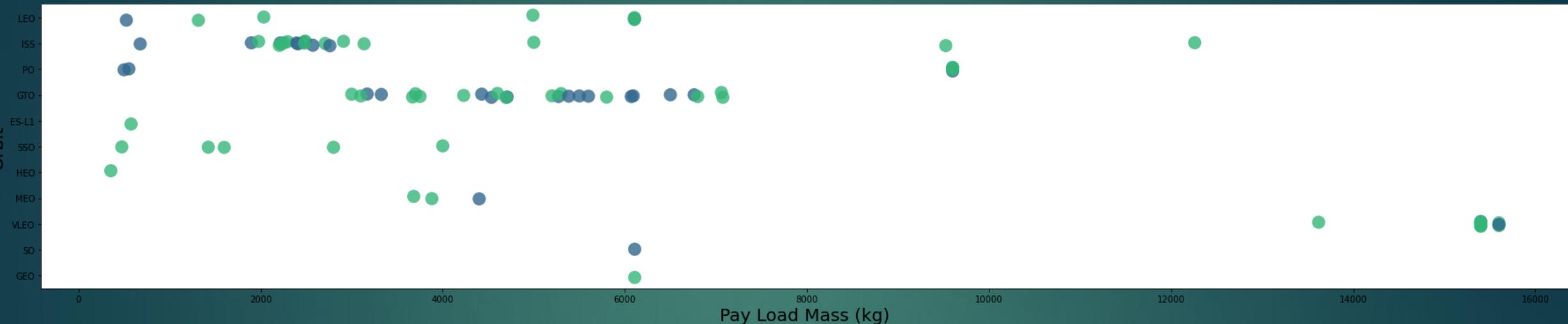
22



Green indicates successful launch; Purple indicates unsuccessful launch.

- Launch Orbit preferences changed over Flight Number.
- Launch Outcome seems to correlate with this preference.
- SpaceX started with LEO orbits which saw moderate success LEO and returned to VLEO in recent launches
- SpaceX appears to perform better in lower orbits or Sun-synchronous orbits

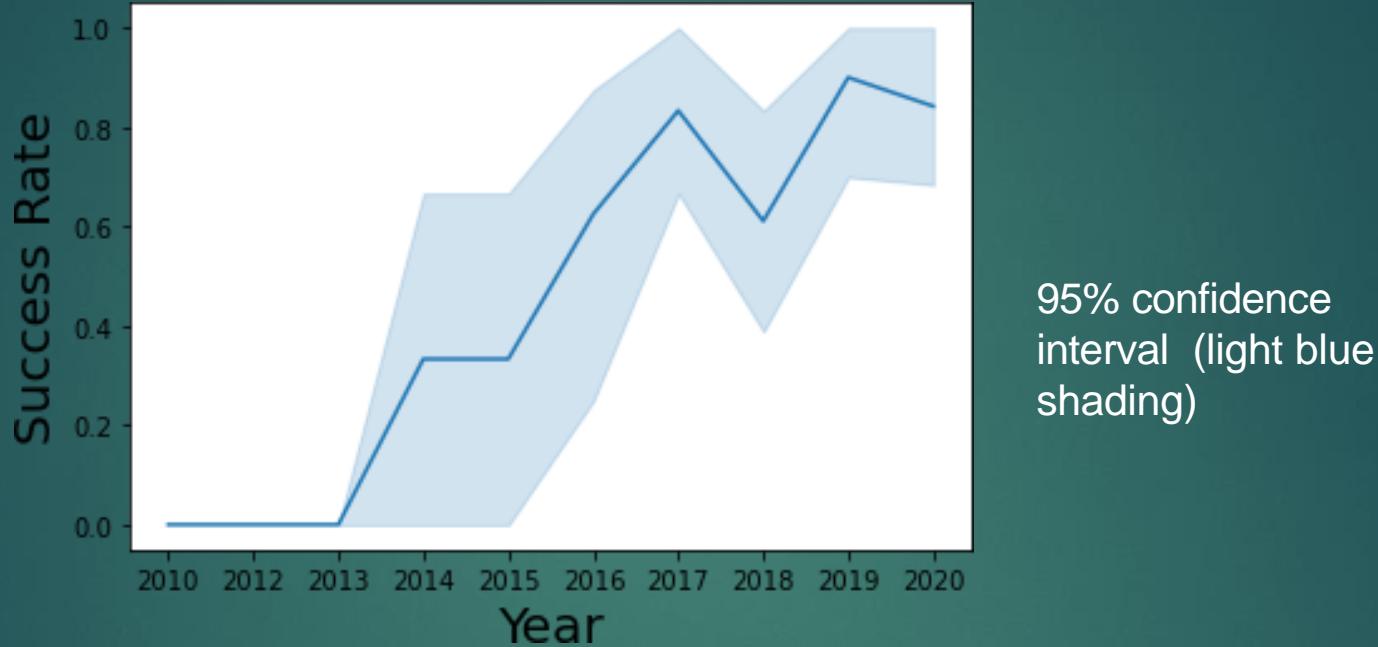
Payload vs. Orbit type



Green indicates successful launch; Purple indicates unsuccessful launch.

- Payload mass seems to correlate with orbit
- LEO and SSO seem to have relatively low payload mass
- The other most successful orbit VLEO only has payload mass values in the higher end of the range

Launch Success Yearly Trend



95% confidence
interval (light blue
shading)

- Success generally increases over time since 2013 with a slight dip in 2018
- Success in recent years at around 80%

EDA - WITH SQL

All Launch Site Names

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

Query unique launch site names from database.

CCAFS SLC-40 and CCAFSSL-40 likely all represent the same launch site with data entry errors.

CCAFS LC-40 was the previous name.

Likely only 3 unique launch_site values:

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



```
1 %sql SELECT UNIQUE(LAUNCH_SITE) FROM SPACEXTBL;
```

Launch Site Names Beginning with `CCA`

```
1 %sql SELECT LAUNCH_SITE FROM SPACEXTBL WHERE LAUNCH_SITE LIKE 'CCA%' LIMIT 5;
```

First five entries in database with Launch Site name beginning with CCA.

launch_site
CCAFS LC-40

Total Payload Mass from NASA

28

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



total_payload_mass
45596

This query sums the total payload mass in kg where NASA was the customer.

CRS stands for Commercial Resupply Services which indicates that these payloads were sent to the International Space Station (ISS).

Average Payload Mass by F9v1.1

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



average_payload_mass
2928

This query calculates the average payload mass or launches which used booster version F9 v1.1

Average payload mass of F9 1.1 is on the low end of our payload mass range

First Successful Ground Pad Landing Date

30

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



first_successful_ground_landing

2015-12-22

This query returns the first successful ground pad landing date.

First ground pad landing wasn't until the end of 2015.

Successful landings in general appear starting 2014.

Successful Drone Ship Landing with Payload Between 4000 and 6000

1

```
%sql SELECT BOOSTER_VERSION FROM SPACEXTBL \
WHERE (LANDING_OUTCOME = 'Success (drone ship)') AND (PAYLOAD_MASS_KG_ BETWEEN 4000 AND 6000);
```



booster_version

F9 FT B1022

F9 FT B1026

F9 FT B1021.2

F9 FT B1031.2

This query returns the four booster versions that had successful drone ship landings and a payload mass between 4000 and 6000 noninclusively.

Total Number of Each Mission Outcome

```
1 %sql SELECT MISSION_OUTCOME, COUNT(MISSION_OUTCOME) AS TOTAL_NUMBER FROM SPACEXTBL GROUP BY MISSION_OUTCOME;
```



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

....s query returns a count of each mission outcome.

SpaceX appears to achieve its mission outcome nearly 99% of the time.

This means that most of the landing failures are intended.

Interestingly, one launch has an unclear payload status and unfortunately one failed in flight.

Boosters that Carried Maximum Payload



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



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

This query returns the booster versions that carried the highest payload mass of 15600 kg.

These booster versions are very similar and all are of the F9 B5 B10xx.x variety.

This likely indicates payload mass correlates with the booster version that is used.

2015 Failed Drone Ship Landing Records



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



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

This query returns the Month, Landing Outcome, Booster Version, Payload Mass (kg), and Launch site of 2015 launches where stage 1 failed to land on a drone ship.

There were two such occurrences.

Ranking Counts of Successful Landings Between 2010-06-04 and 2017-03-20



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



landing_outcome	total_number
No attempt	10
Failure (drone ship)	5
Success (drone ship)	5
Controlled (ocean)	3
Success (ground pad)	3
Failure (parachute)	2
Uncontrolled (ocean)	2
Precluded (drone ship)	1

This query returns a list of successful landings and between 2010-06-04 and 2017-03-20 inclusively.

There are two types of successful landing outcomes: drone ship and ground pad landings.

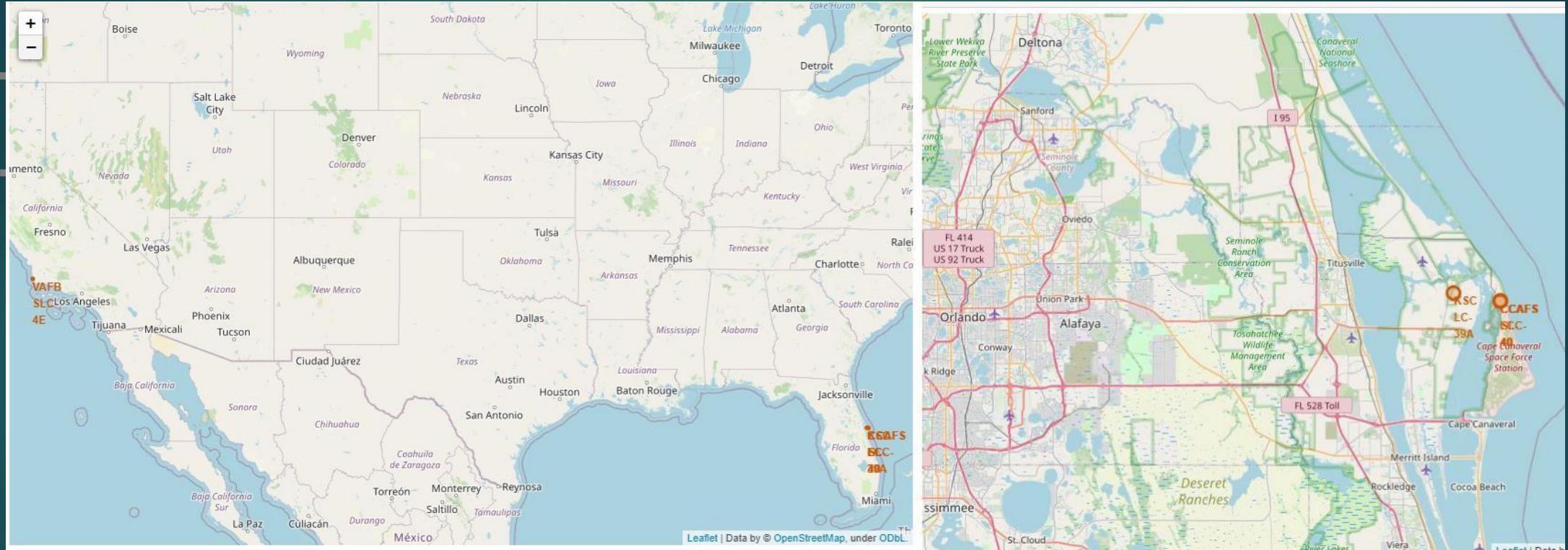
There were 8 successful landings in total during this time period

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. Numerous glowing yellow and white points represent city lights, concentrated in coastal and urban areas. In the upper right quadrant, there are bright green and yellow bands of light, likely the Aurora Borealis or Australis. The overall atmosphere is dark and mysterious.

Section 3

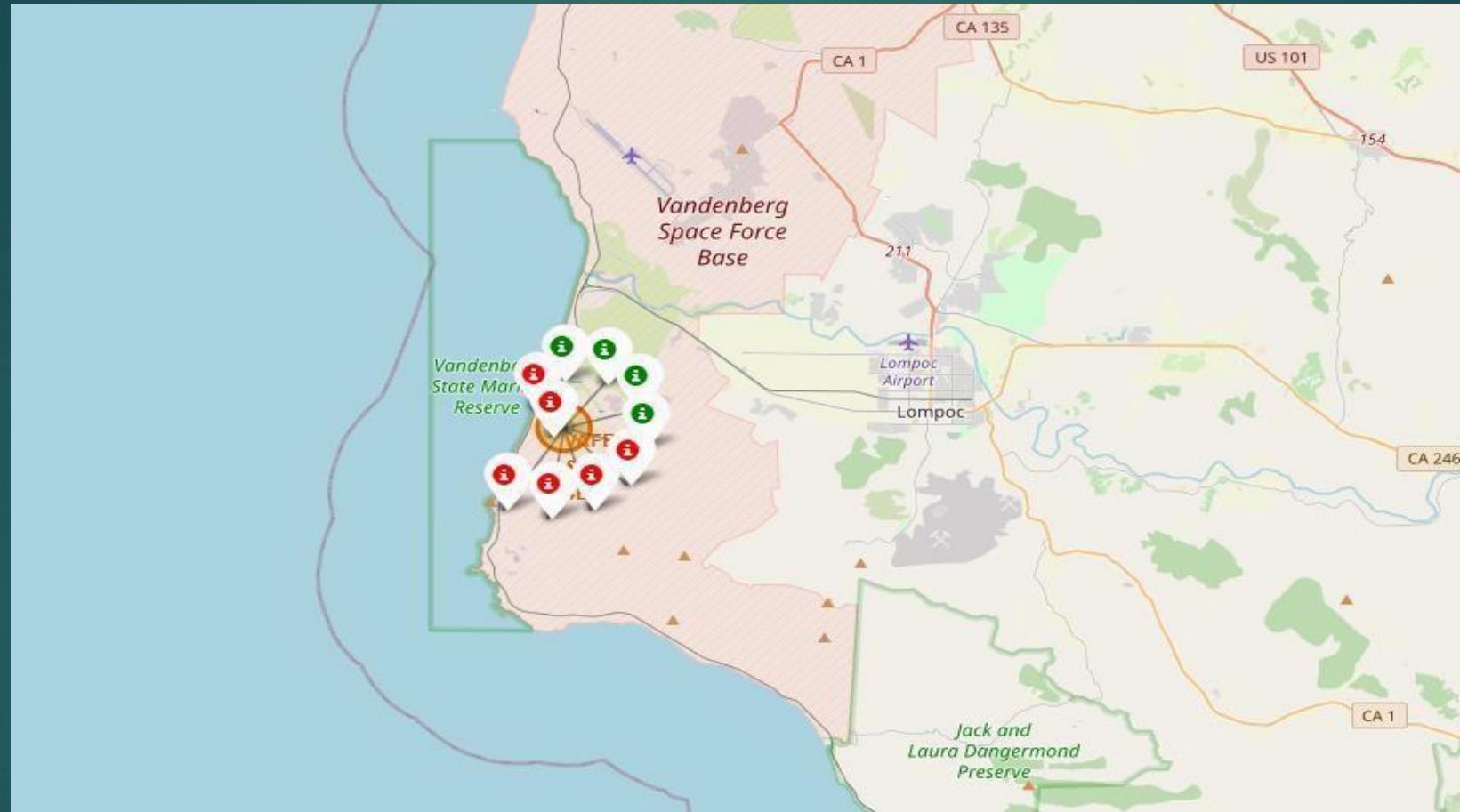
Launch Sites Proximities Analysis

Launch Site



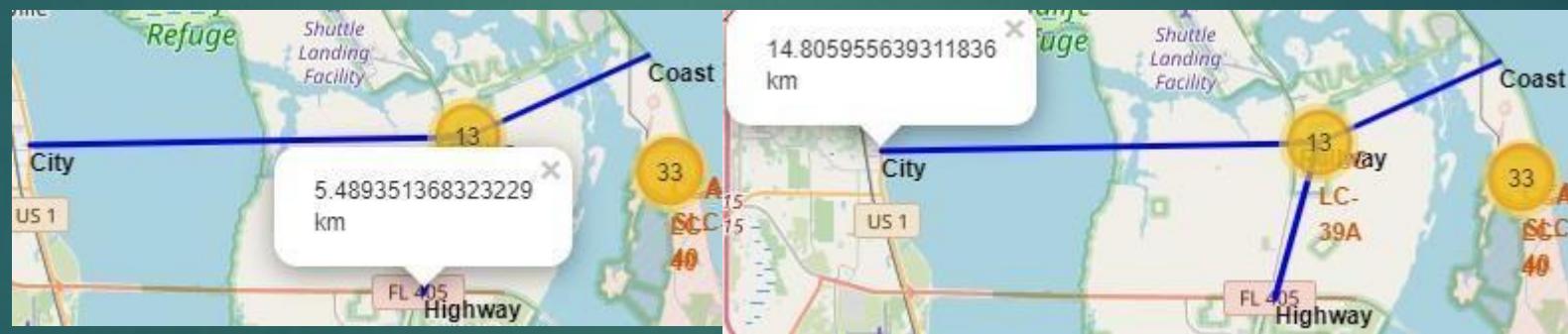
The left map shows all launch sites relative US map. The right map shows the two Florida launch sites since they are very close to each other. All launch sites are near the ocean.

Color-Coded Launch Markers

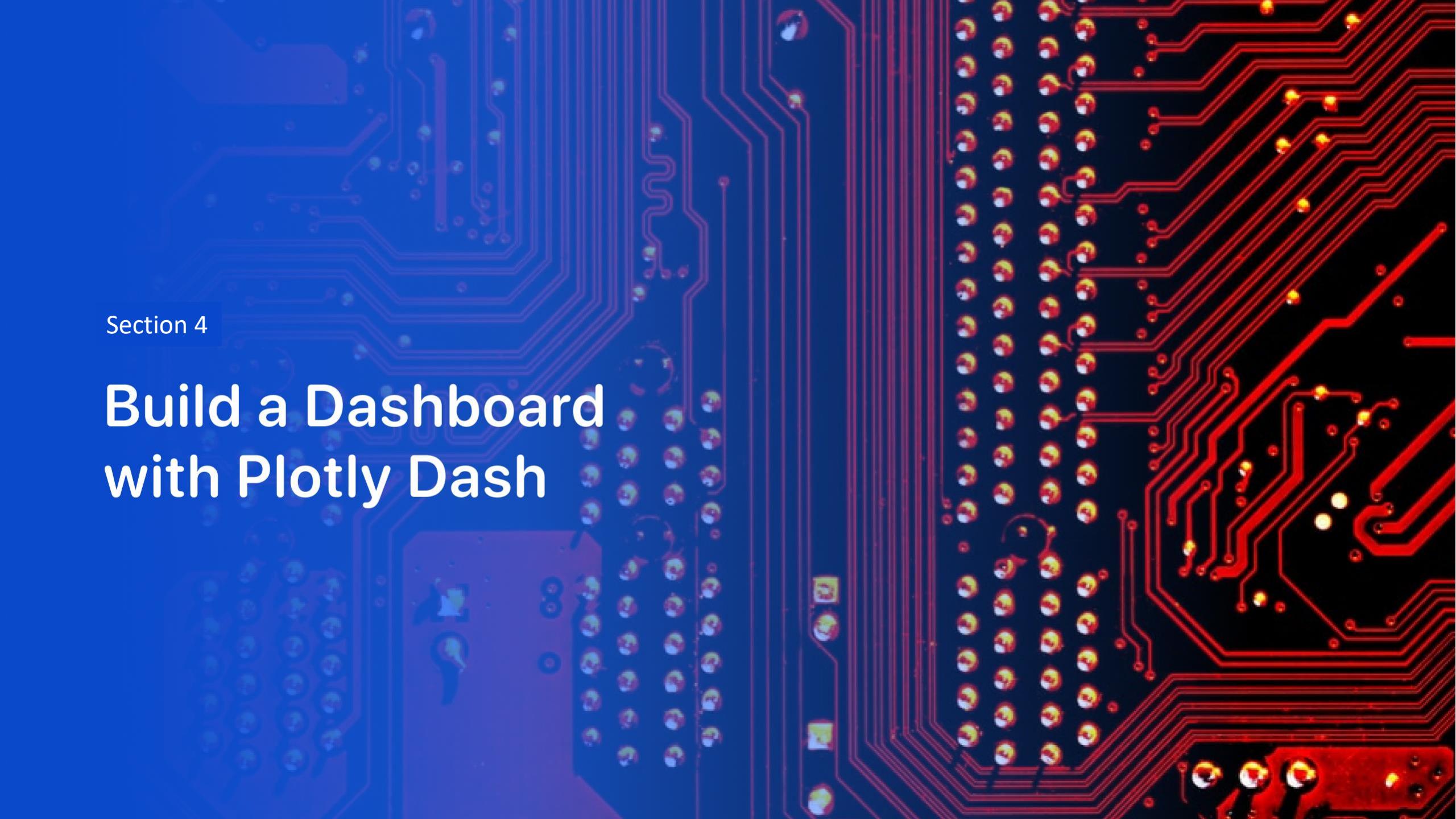


Clusters on Folium map can be clicked on to display each successful landing (green icon) and failed landing (red icon). In this example VAFB SLC-4E shows 4 successful landings and 6 failed landings.

Key Location Proximities



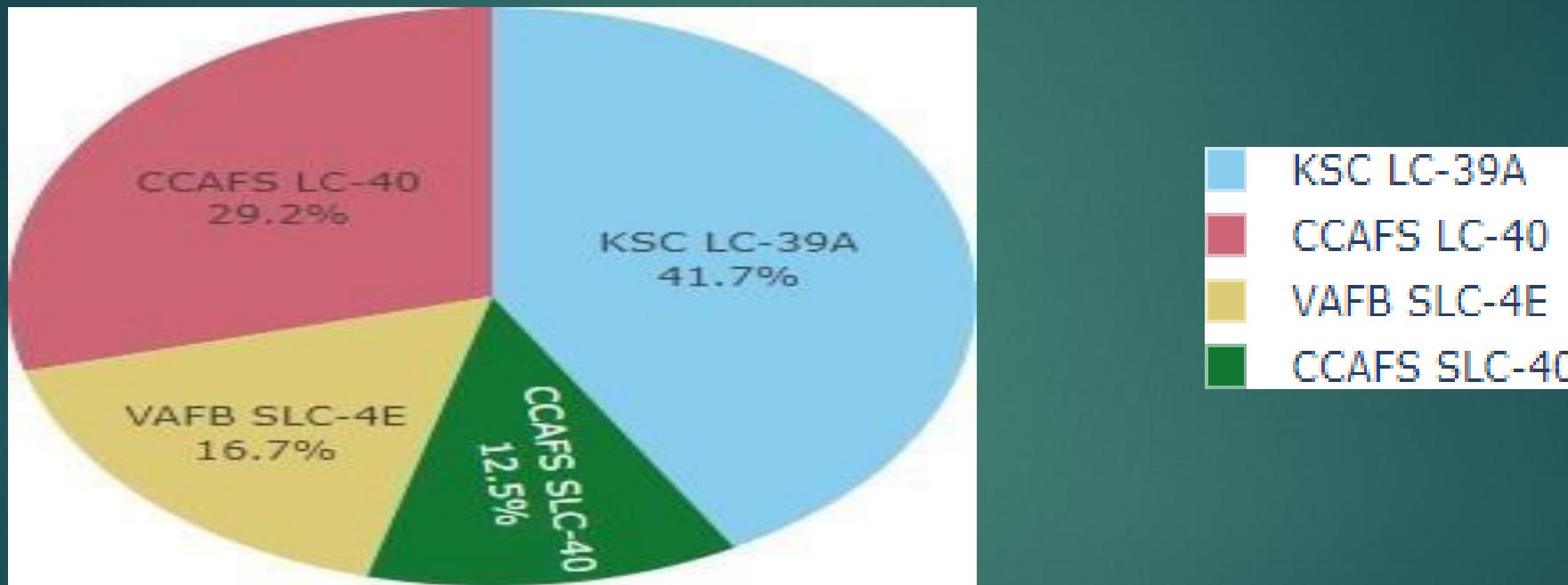
Using KSC LC-39A as an example, launch sites are very close to railways for large part and supply transportation. Launch sites are close to highways for human and supply transport. Launch sites are also close to coasts and relatively far from cities so that launch failures can land in the sea to avoid rockets falling on densely populated areas.



Section 4

Build a Dashboard with Plotly Dash

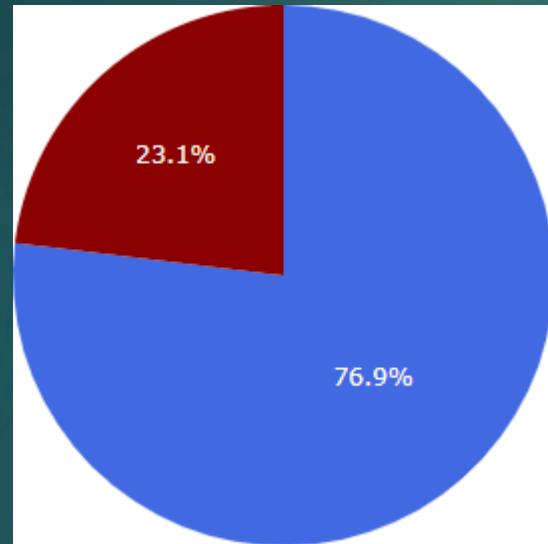
Successful Launches Across Launch Sites



This is the distribution of successful landings across all launch sites. CCAFS LC-40 is the old name of CCAFS SLC-40 so CCAFS and KSC have the same amount of successful landings, but a majority of the successful landings were performed before the name change. VAFB has the smallest share of successful landings. This may be due to smaller sample and increase in difficulty of launching in the west coast.

Highest Success Rate Launch Site

KSC LC-39A Success Rate (blue=success)



KSC LC-39A has the highest success rate with 10 successful landings and 3 failed landings.

Payload Mass vs. Success vs. Booster Version Category



Plotly dashboard has a Payload range selector. However, this is set from 0-10000 instead of the max Payload of 15600. Class indicates 1 for successful landing and 0 for failure. Scatter plot also accounts for booster version category in color and number of launches in point size. In this particular range of 0-6000, interestingly there are two failed landings with payloads of zero kg.

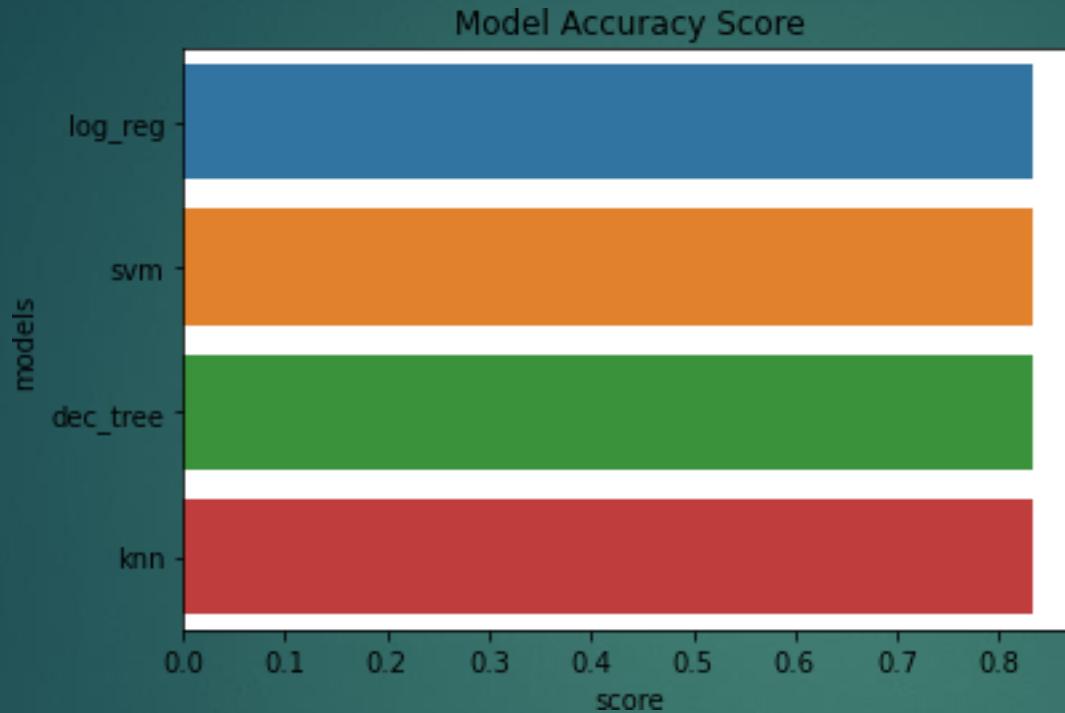
The background of the slide features a dynamic, abstract design. It consists of several thick, curved lines that transition from a bright yellow at the top right to a deep blue at the bottom left. These lines create a sense of motion and depth, resembling a tunnel or a stylized landscape. The overall effect is modern and professional.

Section 5

Predictive Analysis (Classification)

Classification Accuracy

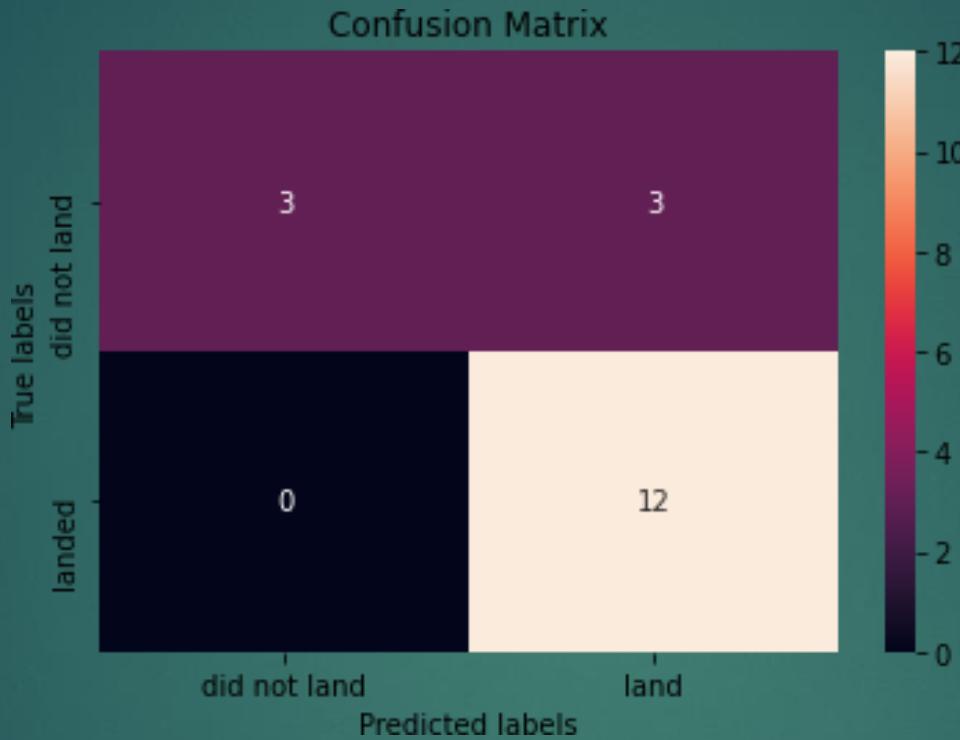
45



- All models had virtually the same accuracy on the test set at 83.33% accuracy.
- It should be noted that test size is small at only sample size of 18.
- This can cause large variance in accuracy results, such as those in Decision Tree Classifier model in repeated runs.
- We likely need more data to determine the best model.

Confusion Matrix

46



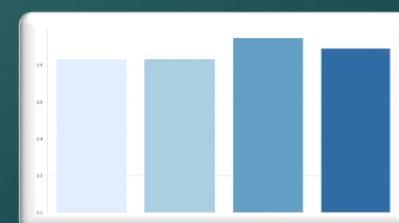
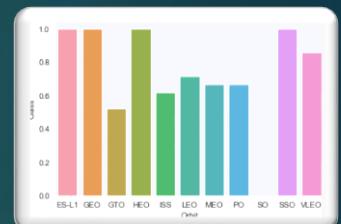
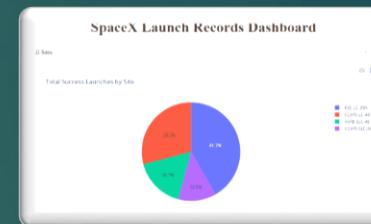
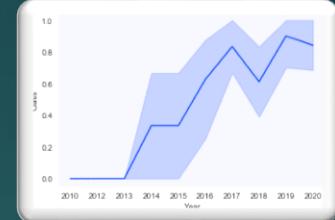
Correct predictions are on a diagonal from top left to bottom right.

- Since all models performed the same for the test set, the confusion matrix is the same across all models.
- The models predicted 12 successful landings when the true label was successful landing.
- The models predicted 3 unsuccessful landings when the true label was unsuccessful landing.
- The models predicted 3 successful landings when the true label was unsuccessful landings (false positives). Our models over predict successful landings.

CONCLUSION

47

- Our task: to develop a machine learning model for Space Y who wants to bid against SpaceX
 - The goal of model is to predict when Stage 1 will successfully land to save ~\$100 million USD
 - Used data from a public SpaceX API and web scraping SpaceX Wikipedia page
 - Created data labels and stored data into a DB2 SQL database
 - Created a dashboard for visualization
 - We created a machine learning model with an accuracy of 83%
 - Allon Mask of SpaceY can use this model to predict with relatively high accuracy whether a launch will have a successful Stage 1 landing before launch to determine whether the launch should be made or not
 - If possible more data should be collected to better determine the best machine learning model and improve accuracy



APPENDIX

48

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

```
def date_time(table_cells):  
    """  
    This function returns the data and time from the HTML table cell  
    Input: the element of a table data cell extracts extra row  
    """  
    return [data_time.strip() for data_time in list(table_cells.strings)][0:2]  
  
def booster_version(table_cells):  
    """  
    This function returns the booster version from the HTML table cell  
    Input: the element of a table data cell extracts extra row  
    """  
    out=''.join([booster_version for i,booster_version in enumerate(table_cells.strings) if i%2==0][0:-1])  
    return out  
  
def landing_status(table_cells):  
    """  
    This function returns the landing status from the HTML table cell  
    Input: the element of a table data cell extracts extra row  
    """  
    out=[i for i in table_cells.strings][0]  
    return out  
  
def get_mass(table_cells):  
    mass=unicodedata.normalize("NFKD", table_cells.text).strip()  
    if mass:  
        mass.find("kg")  
        new_mass=mass[0:mass.find("kg")+2]  
    else:  
        new_mass=0  
    return new_mass  
  
def extract_column_from_header(row):  
    """  
    This function returns the landing status from the HTML table cell  
    Input: the element of a table data cell extracts extra row  
    """  
    if (row.br):  
        row.br.extract()  
    if row.a:  
        row.a.extract()  
    if row.sup:  
        row.sup.extract()  
  
    column_name = ' '.join(row.contents)  
  
    # Filter the digit and empty names  
    if not(column_name.strip().isdigit()):  
        column_name = column_name.strip()  
    return column_name
```

```
extracted_row = 0  
#Extract each table  
for table_index,table in enumerate(soup.find_all('table','wikitable plainrowheaders collapsible')):  
    # get Table Row  
    for rows in table.find_all('tr'): #  
        #check to see if first table heading is as number corresponding to Launch a number  
        if rows.th:  
            if rows.th.string:  
                flight_number=rows.th.string.strip()  
                flag=flight_number.isdigit()  
            else:  
                flag=False  
        #get table element  
        row=rows.find_all('td')  
        #if it is number save cells in a dictionary  
        if flag:  
            extracted_row += 1  
            # Flight Number value  
            # Append the flight_number into launch_dict with key 'Flight No.'  
            launch_dict["Flight No."].append(flight_number)  
  
            # Date value  
            #Append the date into launch_dict with key 'Date'  
            datatimelist=date_time(rows[0])  
            date = datatimelist[0].strip(',')  
            launch_dict["Date"].append(date)  
  
            # Time value  
            #Append the time into launch_dict with key 'Time'  
            time = datatimelist[1]  
            launch_dict["Time"].append(time)  
  
            # Booster version  
            #Append the bv into launch_dict with key 'Version Booster'  
            bv=booster_version(rows[1])  
            if not(bv):  
                bv=rows[1].a.string  
            launch_dict["Version Booster"].append(bv)  
  
            # Launch Site  
            #Append the bv into launch_dict with key 'Launch site'  
            launch_site = rows[2].a.string  
            launch_dict["Launch site"].append(launch_site)  
  
            # Payload  
            #Append the payload into launch_dict with key 'Payload'  
            payload = rows[3].a.string  
            launch_dict["Payload"].append(payload)  
  
            # Payload Mass  
            #Append the payload_mass into launch_dict with key 'Payload mass'  
            payload_mass = get_mass(rows[4])  
            launch_dict["Payload mass"].append(payload_mass)  
  
            # Orbit  
            #Append the orbit into launch_dict with key 'Orbit'  
            orbit = rows[5].a.string  
            launch_dict["Orbit"].append(orbit)  
  
            # Customer  
            #Append the customer into launch_dict with key 'Customer'  
            if rows[6].a != None:  
                customer = rows[6].a.string  
            else:  
                customer = 'None'  
            launch_dict["Customer"].append(customer)  
  
            # Launch outcome  
            #Append the Launch_outcome into launch_dict with key 'Launch outcome'  
            launch_outcome = list(rows[7].strings)[0]  
            launch_dict["Launch outcome"].append(launch_outcome)  
  
            # Booster Landing  
            #Append the booster_landing into launch_dict with key 'Booster Landing'  
            booster_landing = landing_status(rows[8])  
            launch_dict["Booster landing"].append(booster_landing)  
  
    print("Flight Number: {flight_number}, Date: {date}, Time: {time} \n" +  
          "Booster Version: {bv}, Launch Site: {launch_site} \n" +  
          "Payload: {payload}, Orbit: {orbit} \n" +  
          "Customer: {customer}, Launch Outcome: {launch_outcome} \n" +  
          "Booster Landing: {booster_landing} \n" +  
          "*** ")
```

Thank you!

