



IBM Developer  
SKILLS NETWORK

# Winning Space Race with Data Science

Siddharath Narayan Shakya  
16/07/2023



# Outline

---

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

# Executive Summary

---

- Data used in this project is collected from public SpaceX API and SpaceX Wikipedia page. Created labels column 'class' which classifies successful landings. Data is explored 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. Data is standardized and GridSearchCV is used on it to find best parameters for machine learning models and model is verified using accuracy score.
- 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

---

## Background:

- Commercial Space Age is Here
- Space X has best pricing (\$62 million vs. \$165 million USD)
- Largely due to ability to recover part of rocket (Stage 1)
- Space Y wants to compete with Space X

## Problem:

- Space Y tasks us to train a machine learning model to predict successful Stage 1 recovery



Section 1

# Methodology

# Methodology

---

## Executive Summary

- Data collection methodology:
  - Describe how data was collected
- Perform data wrangling
  - Describe how data was processed
- Perform exploratory data analysis (EDA) using visualization and SQL
- Perform interactive visual analytics using Folium and Plotly Dash
- Perform predictive analysis using classification models
  - How to build, tune, evaluate classification models

# Data Collection

---

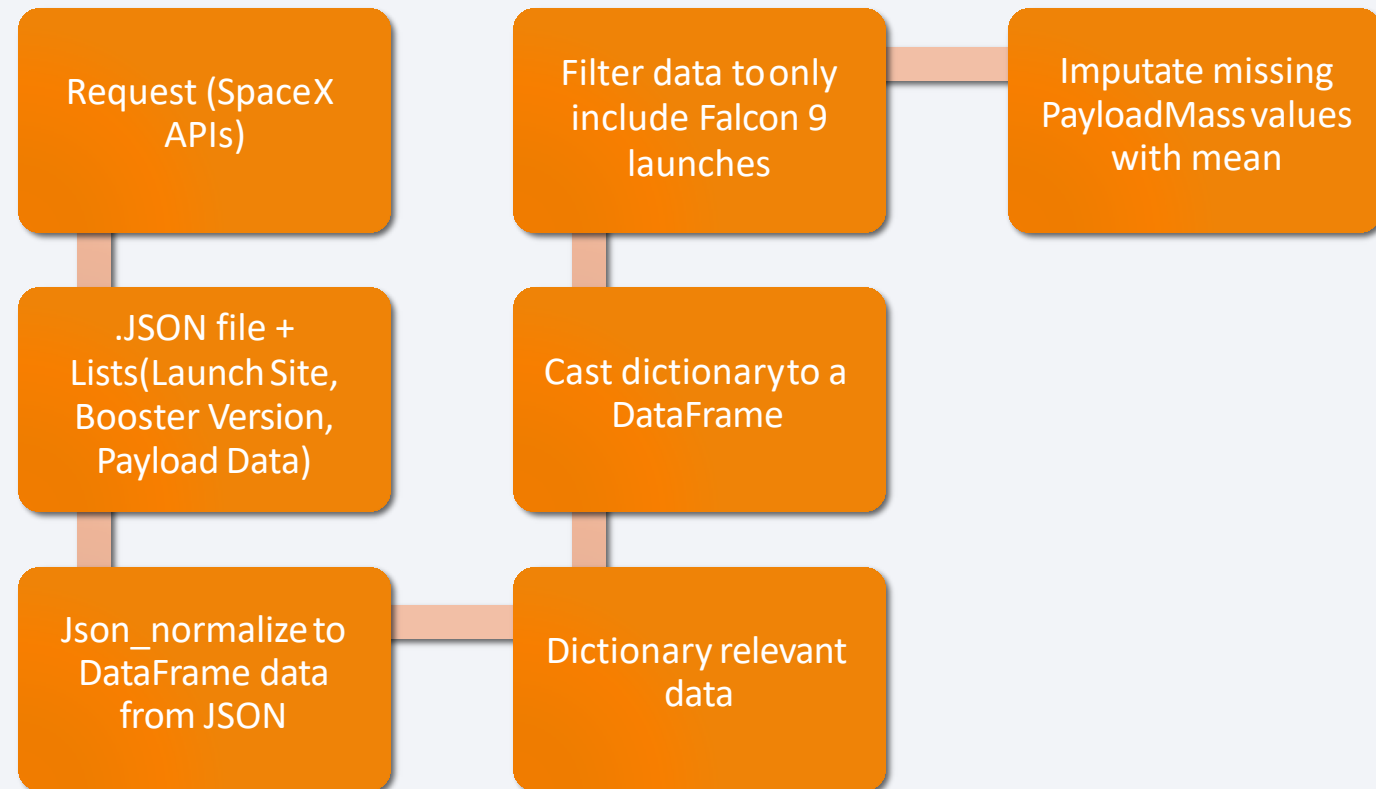
- Data used in this project is collected from public SpaceX API and SpaceX Wikipedia page. Created labels column 'class' which classifies successful landings. Data is explored using SQL, visualization, folium maps, and dashboards. Gathered relevant columns to be used as features. Changed all categorical variables to binary using one hot encoding. Data is standardized and GridSearchCV is used on it to find best parameters for machine learning models and model is verified using accuracy score.
- 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.

# Data Collection – SpaceX API

- SpaceX API Data Columns [FlightNumber, Date, BoosterVersion, PayloadMass, Orbit, LaunchSite, Outcome, Flights, GridFins, Reused, Legs, LandingPad, Block, ReusedCount, Serial, Longitude, Latitude] fetched through REST API Request in the JSON format than converted into Pandas' DataFrame.

- GitHub URL:

- [https://github.com/sidXpro/applied\\_datascience/blob/b2010ecae547e73aa88687a8ce646d4dafa916d1/Applied%20data%20science/week1/jupyter-labs-spacex-data-collection-api.ipynb](https://github.com/sidXpro/applied_datascience/blob/b2010ecae547e73aa88687a8ce646d4dafa916d1/Applied%20data%20science/week1/jupyter-labs-spacex-data-collection-api.ipynb)



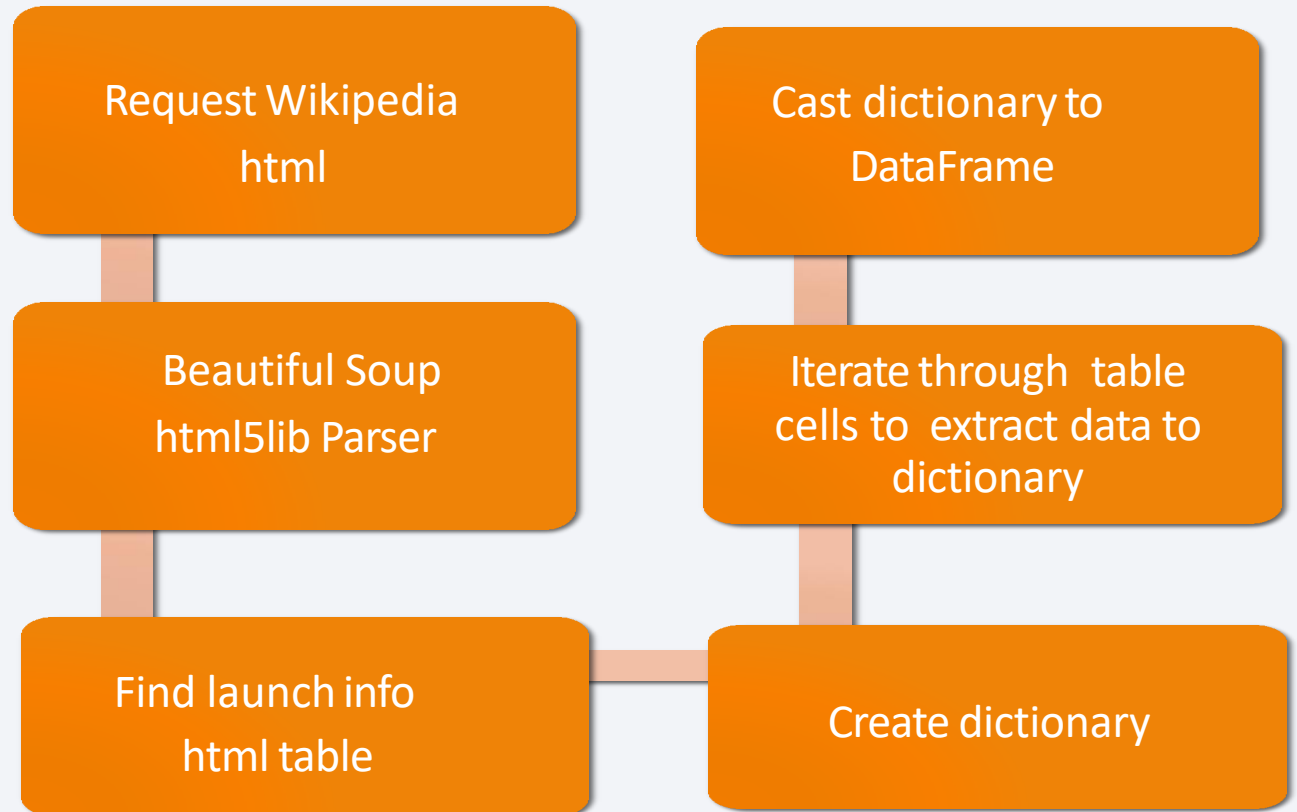


# Data Collection - Scraping

- Wikipedia Webscrape Data Columns includes Flight No., Launch site, Payload, Payload Mass, Orbit, Customer, Launch outcome, Version Booster, Booster landing, Date, Time fetched in html format than the information is extracted from the tables with the help of beautiful soup html parser which is then stored into pandas DataFrame.

- GitHub URL:

[https://github.com/sidXpro/applied\\_datascience/blob/b2010ecae547e73aa88687a8ce646d4dafa916d1/Applied%20data%20science/week1/jupyter-labs-webscraping.ipynb](https://github.com/sidXpro/applied_datascience/blob/b2010ecae547e73aa88687a8ce646d4dafa916d1/Applied%20data%20science/week1/jupyter-labs-webscraping.ipynb)



# 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
- GitHuburl:  
[https://github.com/sidXpro/applied\\_datascience/blob/b2010ecae547e73aa88687a8ce646d4dafa916d1/Applied%20data%20science/week1/SpaceXDataWrangling.ipynb](https://github.com/sidXpro/applied_datascience/blob/b2010ecae547e73aa88687a8ce646d4dafa916d1/Applied%20data%20science/week1/SpaceXDataWrangling.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/sidXpro/applied\\_datascience/blob/b2010ecae547e73aa88687a8ce646d4dafa916d1/Applied%20data%20science/week2/eda\\_datavisualization.ipynb](https://github.com/sidXpro/applied_datascience/blob/b2010ecae547e73aa88687a8ce646d4dafa916d1/Applied%20data%20science/week2/eda_datavisualization.ipynb)

# EDA with SQL

---

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/sidXpro/applied\\_datascience/blob/b2010ecae547e73aa88687a8ce646d4dafa916d1/Applied%20data%20science/week2/jupyter-labs-eda-sql-coursera\\_sqlite.ipynb](https://github.com/sidXpro/applied_datascience/blob/b2010ecae547e73aa88687a8ce646d4dafa916d1/Applied%20data%20science/week2/jupyter-labs-eda-sql-coursera_sqlite.ipynb)

# Build an Interactive Map with Folium

---

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/sidXpro/applied\\_datascience/blob/b2010ecae547e73aa88687a8ce646d4dafa916d1/Applied%20data%20science/week3/launch\\_site\\_location.ipynb](https://github.com/sidXpro/applied_datascience/blob/b2010ecae547e73aa88687a8ce646d4dafa916d1/Applied%20data%20science/week3/launch_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.

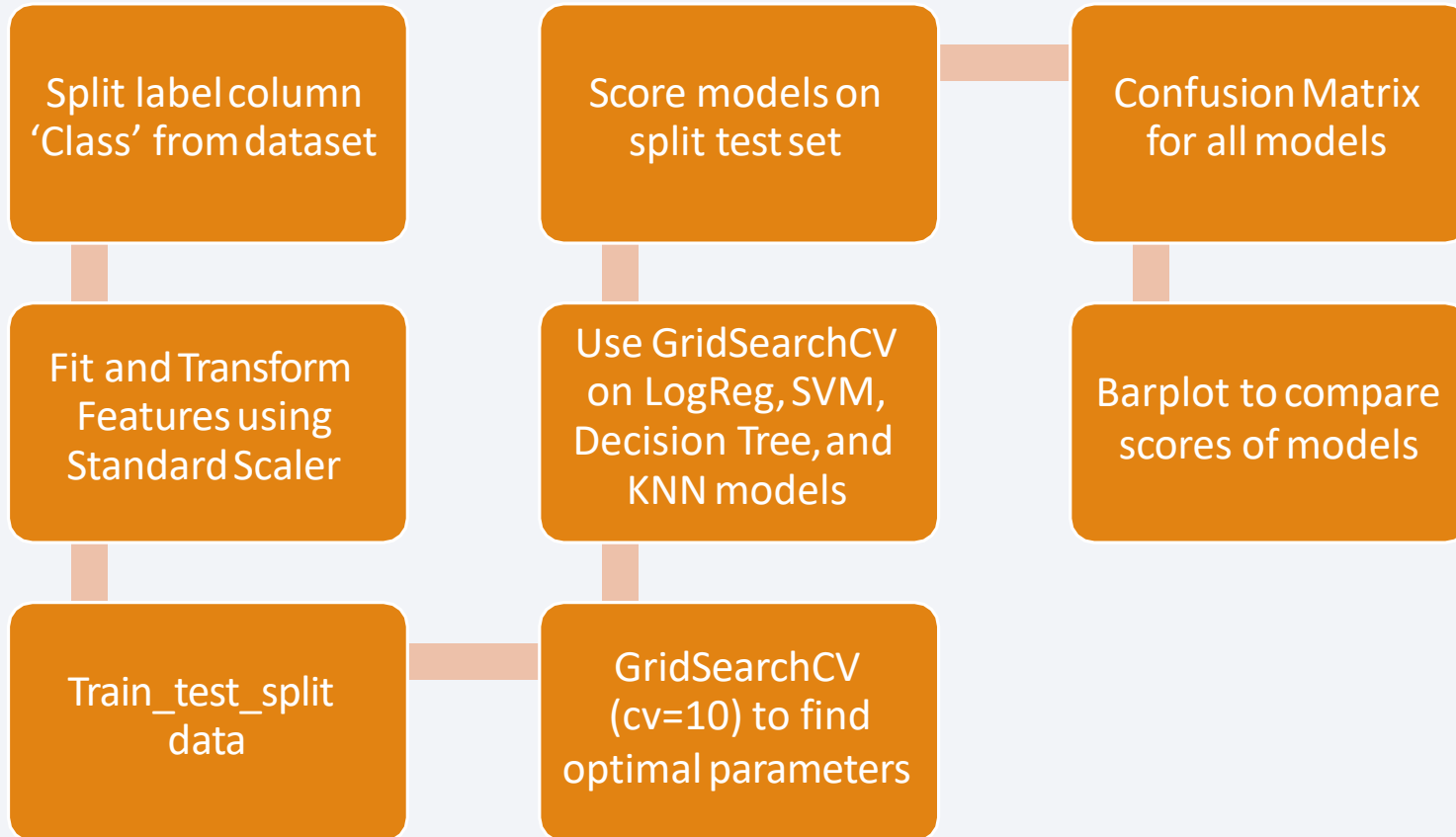
The scatter plot can help us see how success varies across launch sites, payload mass, and booster version category.

GitHub url:

[https://github.com/sidXpro/applied\\_datascience/blob/b2010ecae547e73aa88687a8ce646d4dafa916d1/Applied%20data%20science/week3/spacex\\_dash\\_app.py](https://github.com/sidXpro/applied_datascience/blob/b2010ecae547e73aa88687a8ce646d4dafa916d1/Applied%20data%20science/week3/spacex_dash_app.py)

# Predictive Analysis (Classification)

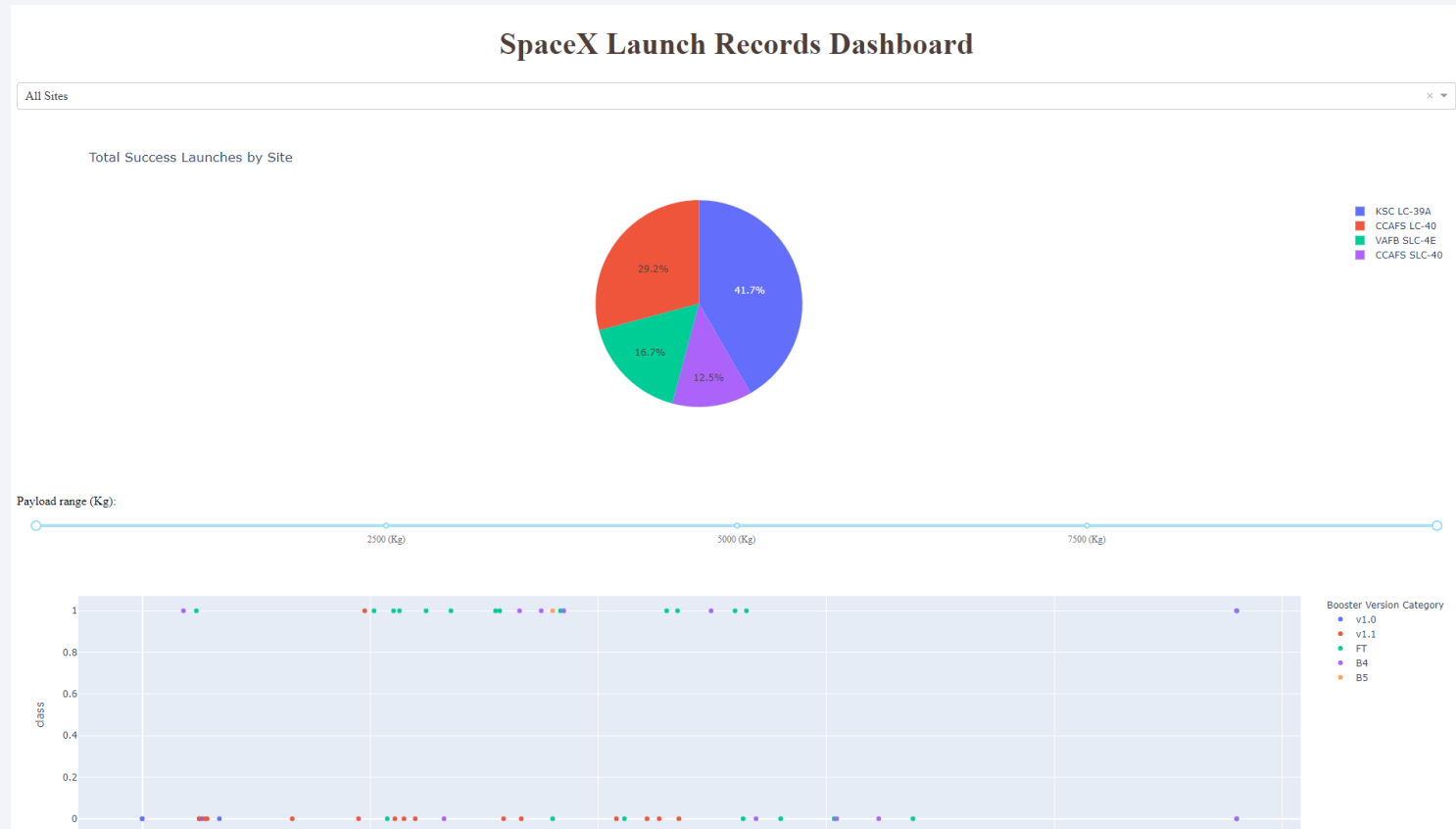
---



GitHub url:

[https://github.com/sidXpro/applied\\_datascience/blob/b2010ecae547e73aa88687a8ce646d4dafa916d1/Applied%20data%20science/week4/MachineLearningPrediction.ipynb](https://github.com/sidXpro/applied_datascience/blob/b2010ecae547e73aa88687a8ce646d4dafa916d1/Applied%20data%20science/week4/MachineLearningPrediction.ipynb)

# Results



This is a preview of the Plotly dashboard. The following slides will show the results of EDA with visualization, EDA with SQL, Interactive Map with Folium, and finally the results of our model with about 83% accuracy.



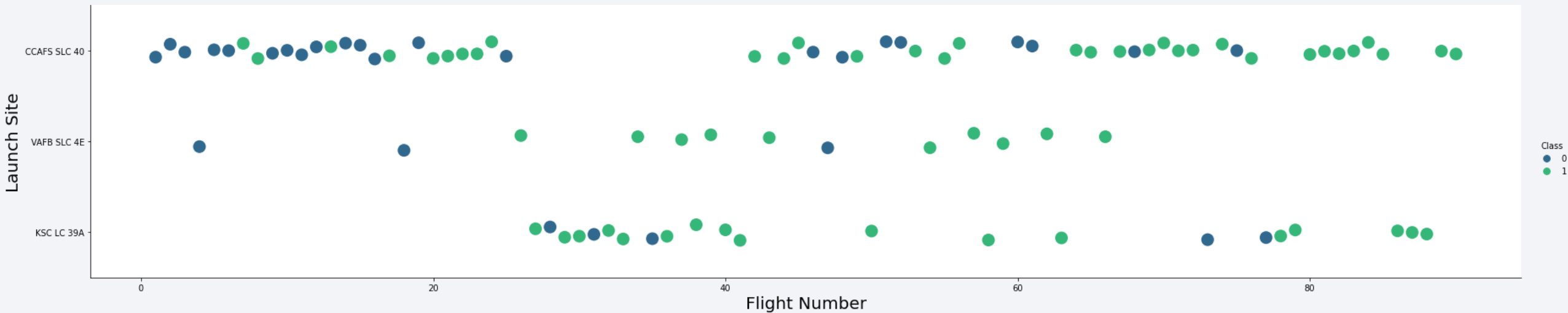
The background of the slide is an abstract composition. It features a dark blue base color. Overlaid on this are numerous diagonal streaks in shades of red and cyan. A faint, light blue grid pattern is also visible, particularly in the lower-left quadrant. The overall effect is dynamic and technological.

Section 2

# Insights drawn from EDA



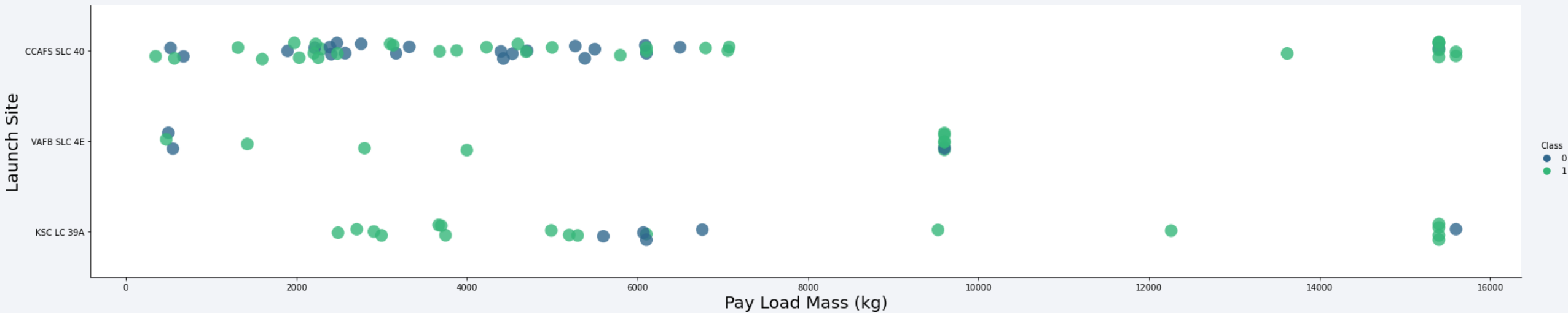
# Flight Number vs. Launch Site



- 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.

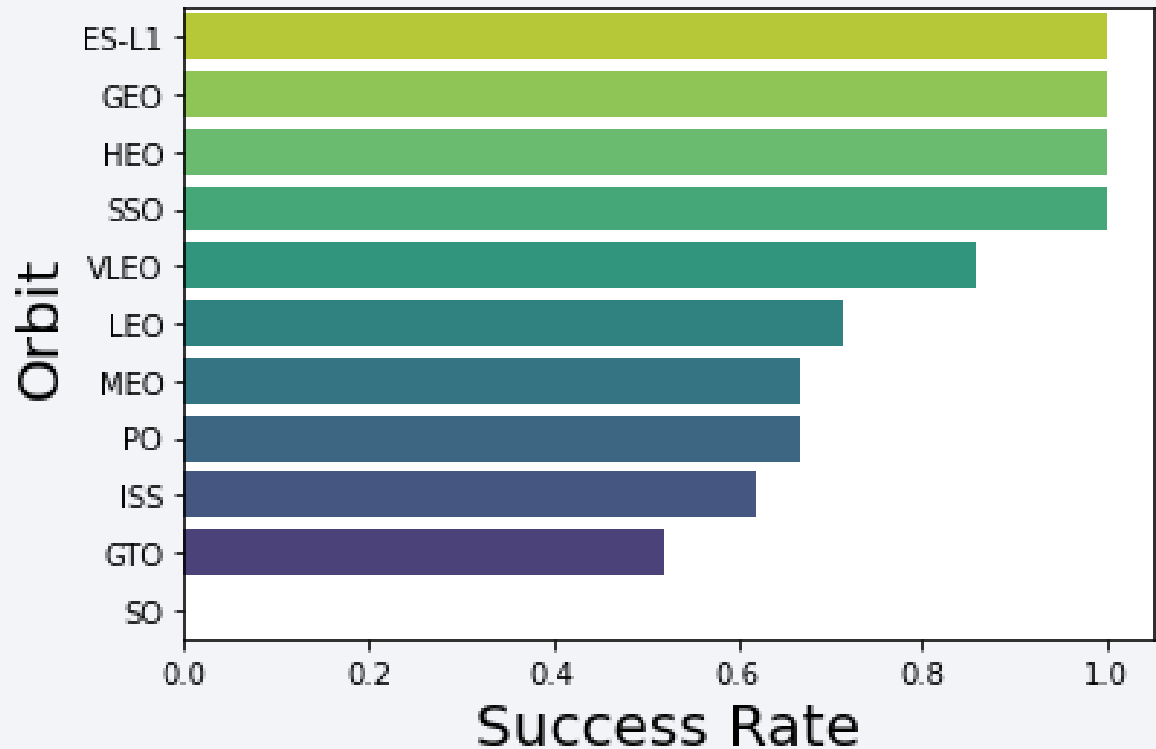


# Payload vs. Launch Site



- 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.

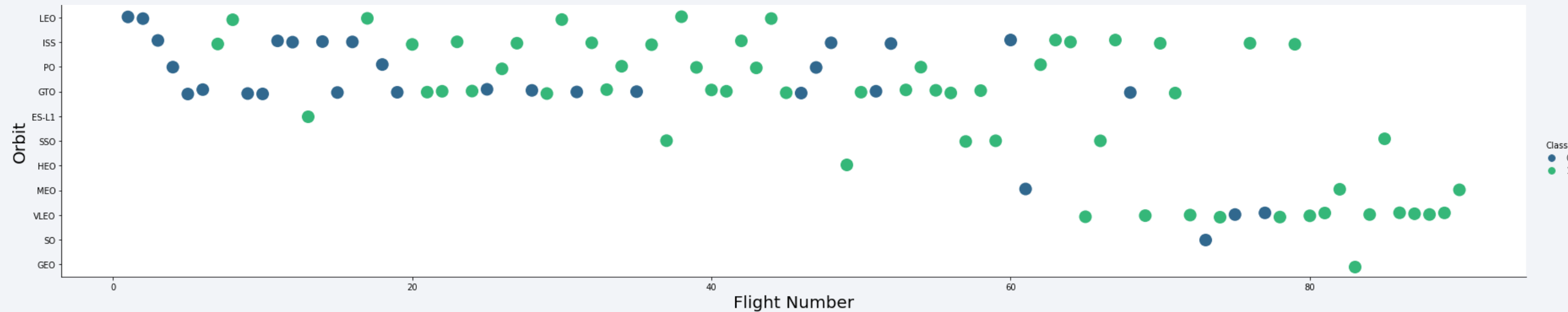
# Success Rate vs. Orbit Type



Success Rate Scale with  
0 as 0%  
0.6 as 60%  
1 as 100%

- 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. Orbit Type



- 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

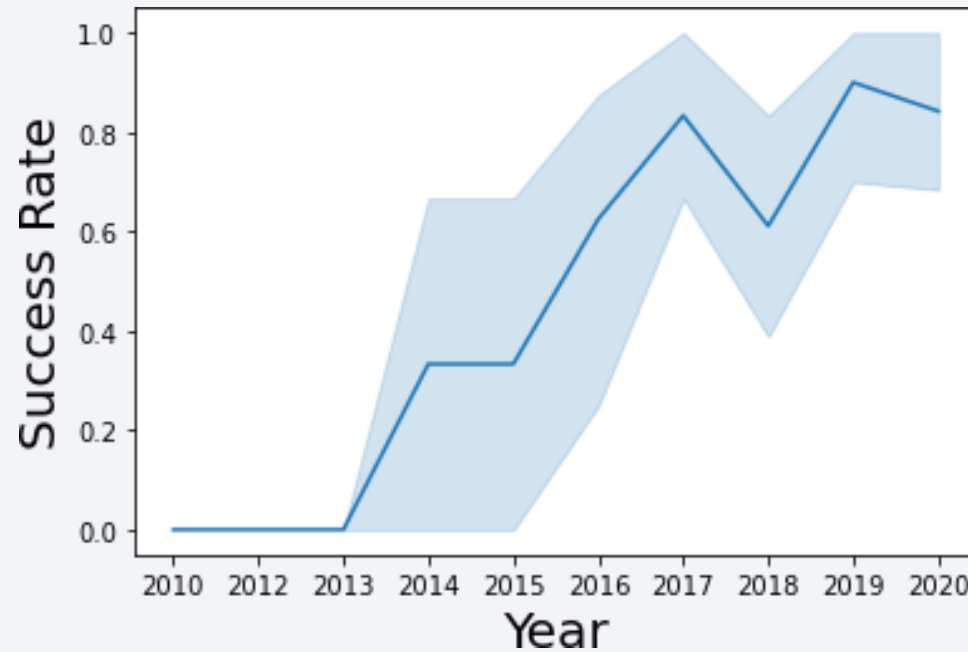
# Payload vs. Orbit Type



- 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%



# All Launch Site Names

---

## Task 1

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

```
In [8]: %sql select DISTINCT LAUNCH_SITE from SPACEXTBL
```

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

```
Out[8]: Launch_Site
```

```
CCAFS LC-40
```

```
VAFB SLC-4E
```

```
KSC LC-39A
```

```
CCAFS SLC-40
```

```
None
```

Query unique launch site names from database.

CCAFS SLC-40 and CCAFSSLC-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

# Launch Site Names Begin with 'CCA'

## Task 2

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

```
In [10]: %sql select * from SPACEXTBL where launch_site like 'CCA%' limit 5
```

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

```
Out[10]:
```

	Date	Time (UTC)	Booster_Version	Launch_Site	Payload	PAYLOAD_MASS_KG_	Orbit	Customer	Mission_Outcome	Landing_Outc
	06/04/2010	18:45:00	F9 v1.0 B0003	CCAFS LC-40	Dragon Spacecraft Qualification Unit	0.0	LEO	SpaceX	Success	Failure (parachute)
	12/08/2010	15:43:00	F9 v1.0 B0004	CCAFS LC-40	Dragon demo flight C1, two CubeSats, barrel of Brouere cheese	0.0	LEO (ISS)	NASA (COTS) NRO	Success	Failure (parachute)
	22/05/2012	7:44:00	F9 v1.0 B0005	CCAFS LC-40	Dragon demo flight C2	525.0	LEO (ISS)	NASA (COTS)	Success	No attachment
	10/08/2012	0:35:00	F9 v1.0 B0006	CCAFS LC-40	SpaceX CRS-1	500.0	LEO (ISS)	NASA (CRS)	Success	No attachment
	03/01/2013	15:10:00	F9 v1.0 B0007	CCAFS LC-40	SpaceX CRS-2	677.0	LEO (ISS)	NASA (CRS)	Success	No attachment

- Sql query to show the launch Site whose names begin with 'CCA' using 'CCA%'
- wild card character(%)

# Total Payload Mass

---

## Task 3

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

```
In [11]: %sql select sum(payload_mass__kg_) as sum from SPACEXTBL where customer like 'NASA (CRS)'
```

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

```
Done.
```

```
Out[11]:
```

sum
45596.0

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 F9 v1.1

---

## Task 4

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

```
In [12]: %sql select avg(payload_mass__kg_) as Average from SPACEXTBL where booster_version like 'F9 v1.1%'
```

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

```
Out[12]:
```

<u>Average</u>
2534.6666666666665

This query calculates the average payload mass of 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 Landing Date

---

## Task 5

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

*Hint: Use min function*

```
In [13]: %sql select min(date) as Date from SPACEXTBL where mission_outcome like 'Success'
```

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

```
Done.
```

```
Out[13]:
```

Date
01/06/2014

This query returns the first successful ground pad landing date.

Successful landings in general appear starting 2014.



# Successful Drone Ship Landing with Payload between 4000 and 6000

## Task 6

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

In [16]:

```
%%sql
select booster_version from SPACEXTBL where (mission_outcome like 'Success')
AND (payload_mass__kg_ BETWEEN 4000 AND 6000) AND (landing_outcome like 'Success (drone ship)')
```

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

Done.

Out[16]:

**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 Successful and Failure Mission Outcomes

## Task 7

List the total number of successful and failure mission outcomes

```
In [18]: %sql SELECT mission_outcome, count(*) as Count FROM SPACEXTBL GROUP by mission_outcome ORDER BY mission_outcome
```

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

Done.

```
Out[18]:
```

Mission_Outcome	Count
None	898
Failure (in flight)	1
Success	98
Success	1
Success (payload status unclear)	1

SpaceX appears to achieve its mission outcome nearly 99% of the time. This means that most of the landing failure are intended. Interestingly, one launch has an unclear payload status and unfortunately one failed in flight.

# Boosters Carried Maximum Payload

## Task 8

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

```
In [20]: maxm = %sql select max(payload_mass__kg_) from SPACEXTBL
maxv = maxm[0][0]
%sql select booster_version from SPACEXTBL where payload_mass__kg_=(select max(payload_mass__kg_) from SPACEXTBL)

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

Out[20]: **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

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 Launch Records

## Task 9

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, 4, 2) as month to get the months and substr(Date,7,4)='2015' for year.**

```
In [30]: %%sql
select substr(Date, 4, 2) as month, landing_outcome, booster_version, launch_site
from SPACEXTBL where landing_outcome like 'Failure (drone ship)' AND substr(Date,7,4)='2015'
```

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

```
Done.
```

```
Out[30]:
```

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

There were two such occurrences.

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

---

```
%%sql
SELECT landing__outcome, COUNT(*) AS no_outcome
FROM SPACEXDATASET
WHERE landing__outcome LIKE 'Succes%' AND DATE BETWEEN '2010-06-04' AND '2017-03-20'
GROUP BY landing__outcome
ORDER BY no_outcome DESC;
```

landing__outcome	no_outcome
Success (drone ship)	5
Success (ground pad)	3

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

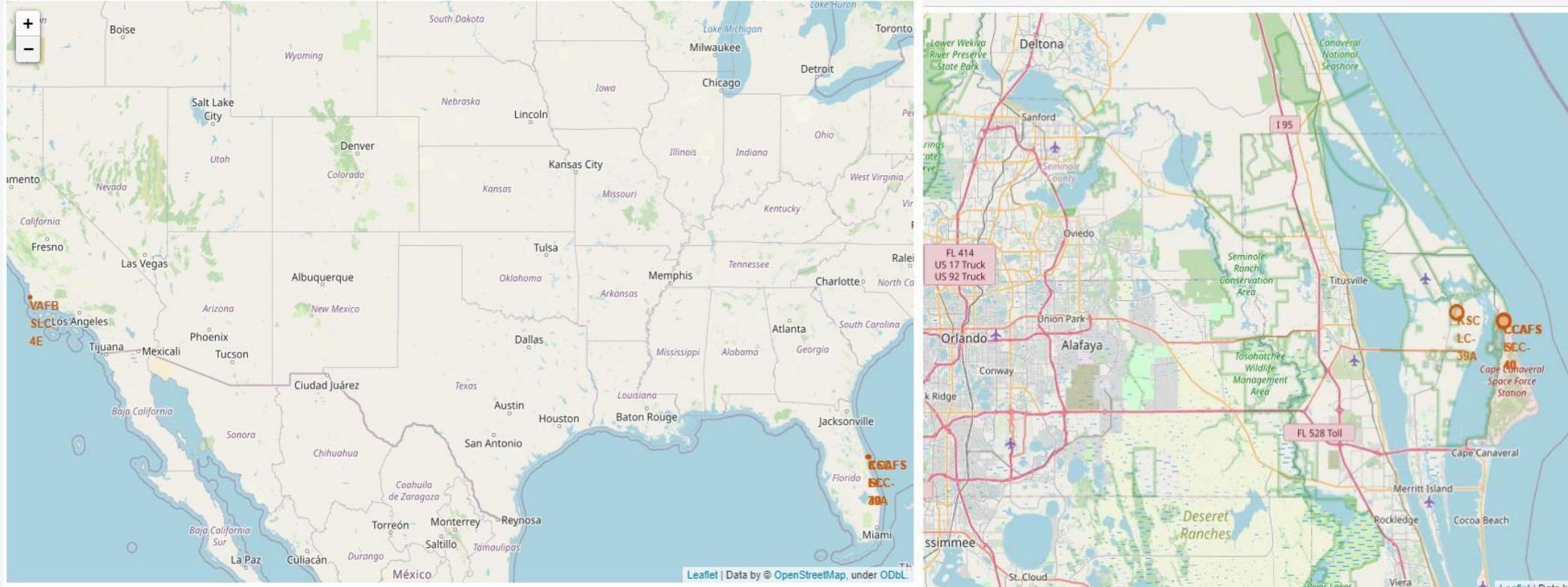
A satellite view of Earth from space, showing the curvature of the planet and city lights at night. The background is a deep blue gradient.

Section 3

# Launch Sites Proximities Analysis

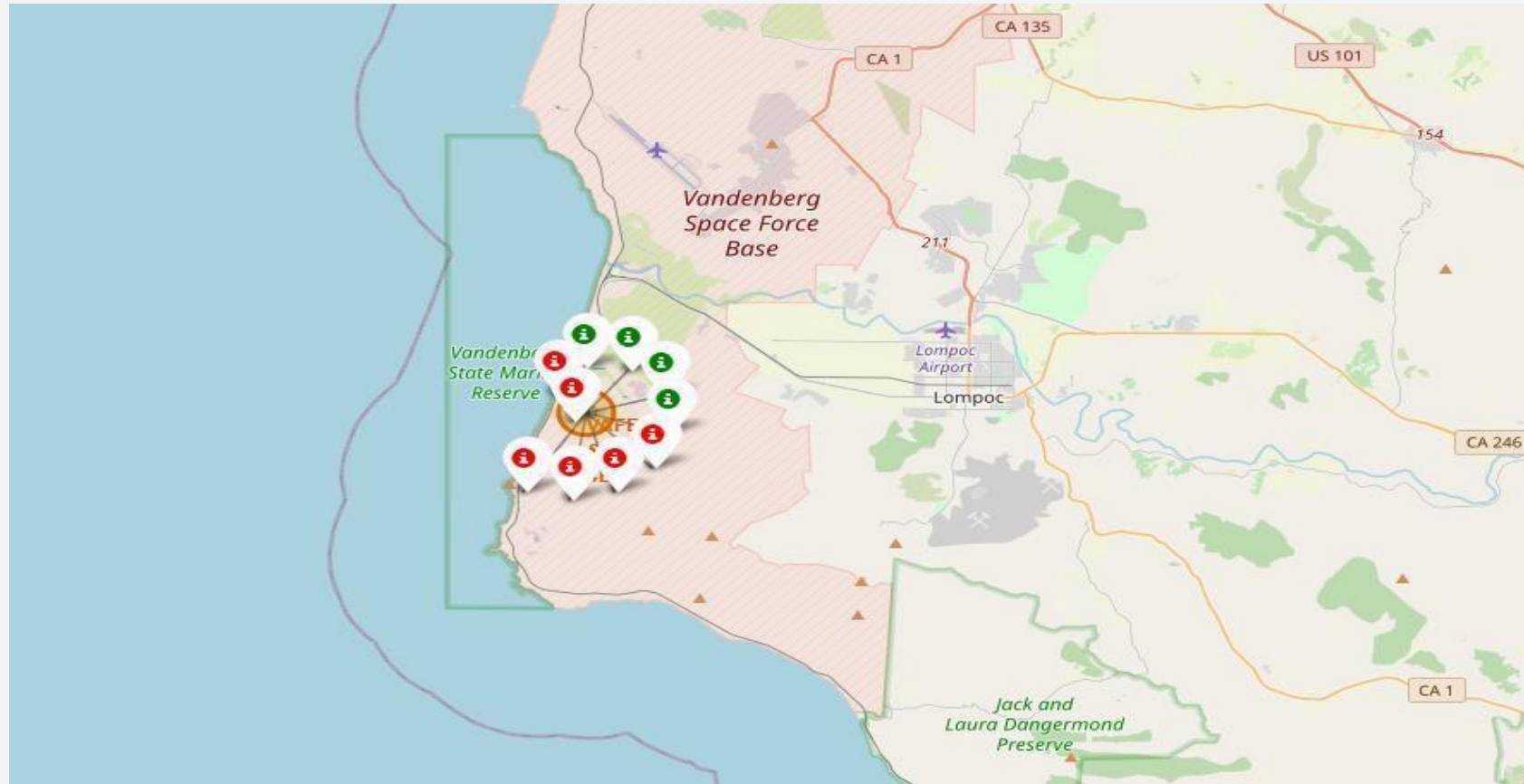


# Launch Site Locations



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.

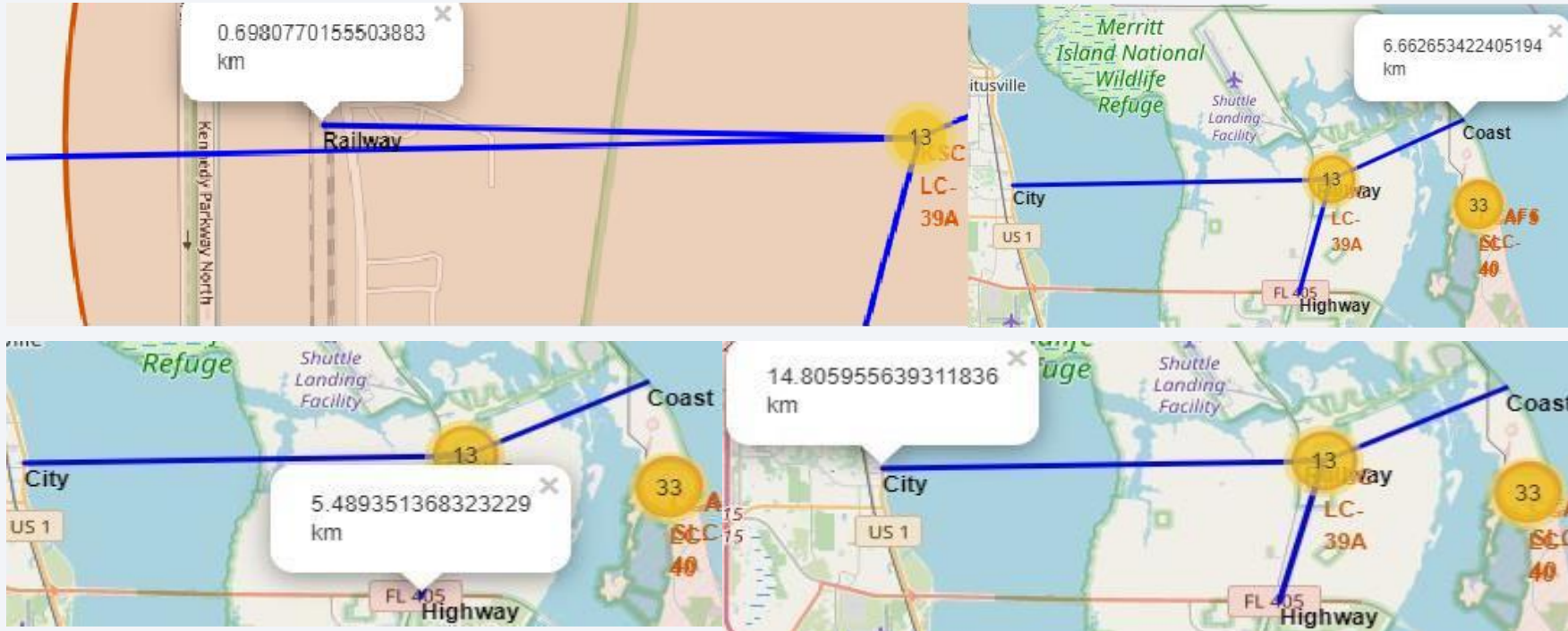
# Color-Coded Launch Markers



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.

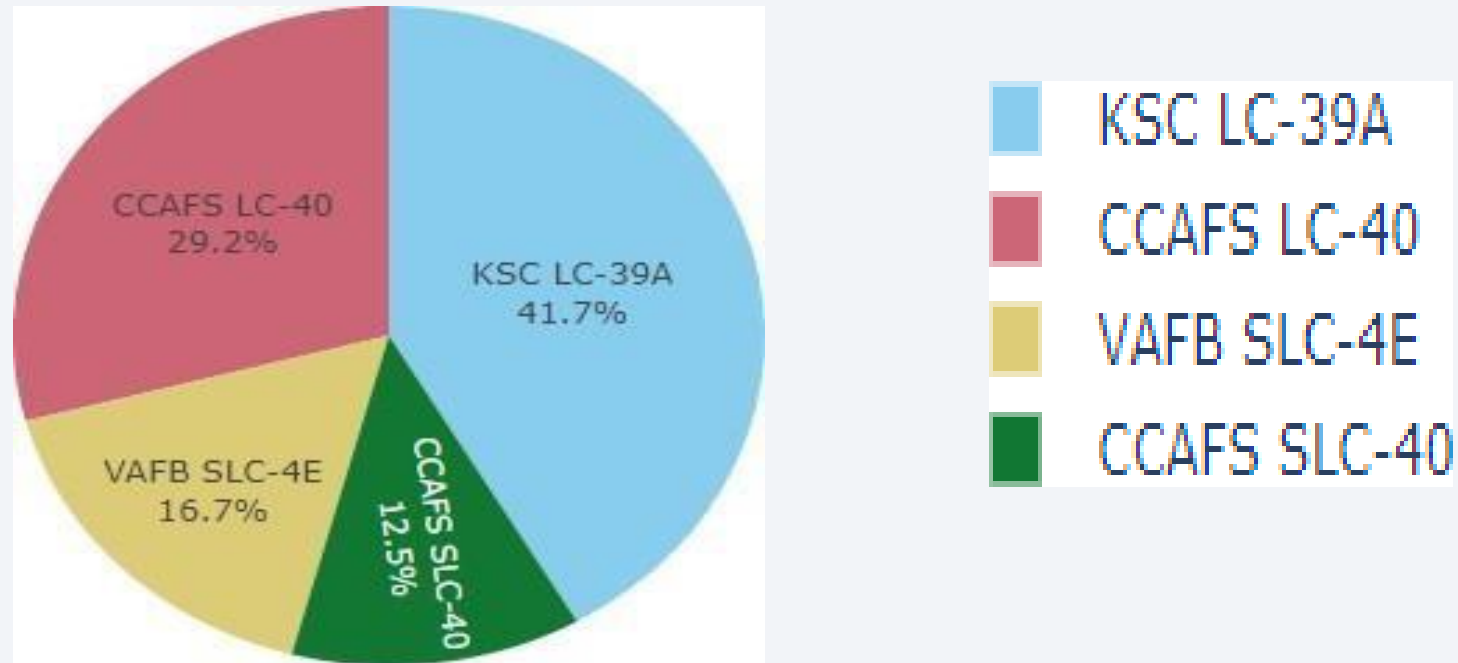


Section 4

# Build a Dashboard with Plotly Dash

# Successful Launches Across Launch Sites

---

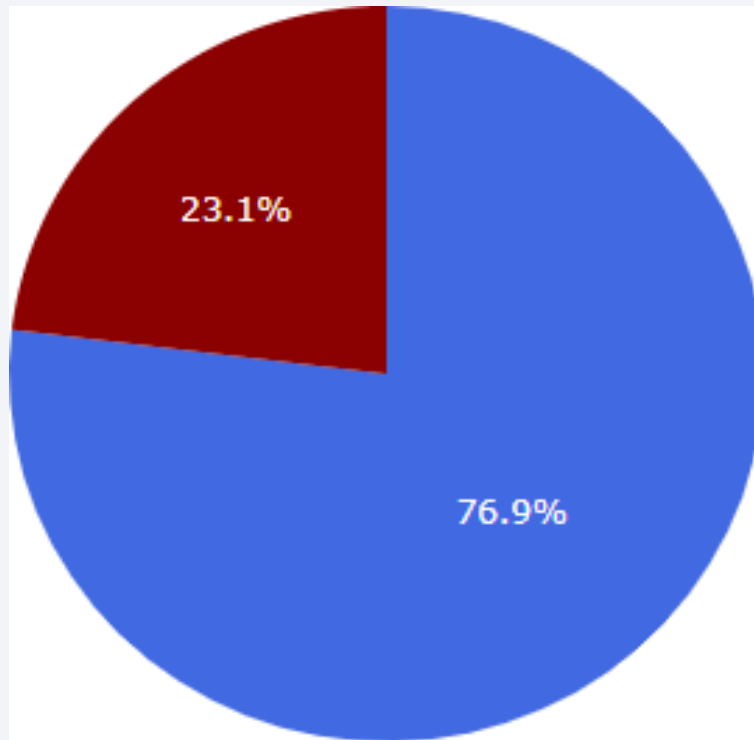


This is the distribution of successful landings across all launch sites. CCAFS LC-40 is the old name of CCAFS SLC-40 so CCAFS and KSC have the same amount of successful landings, but a majority of the successful landings were performed before the name change. VAFB has the smallest share of successful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

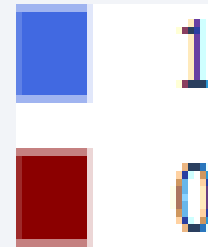


# High Success rate Launch Sites

---



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



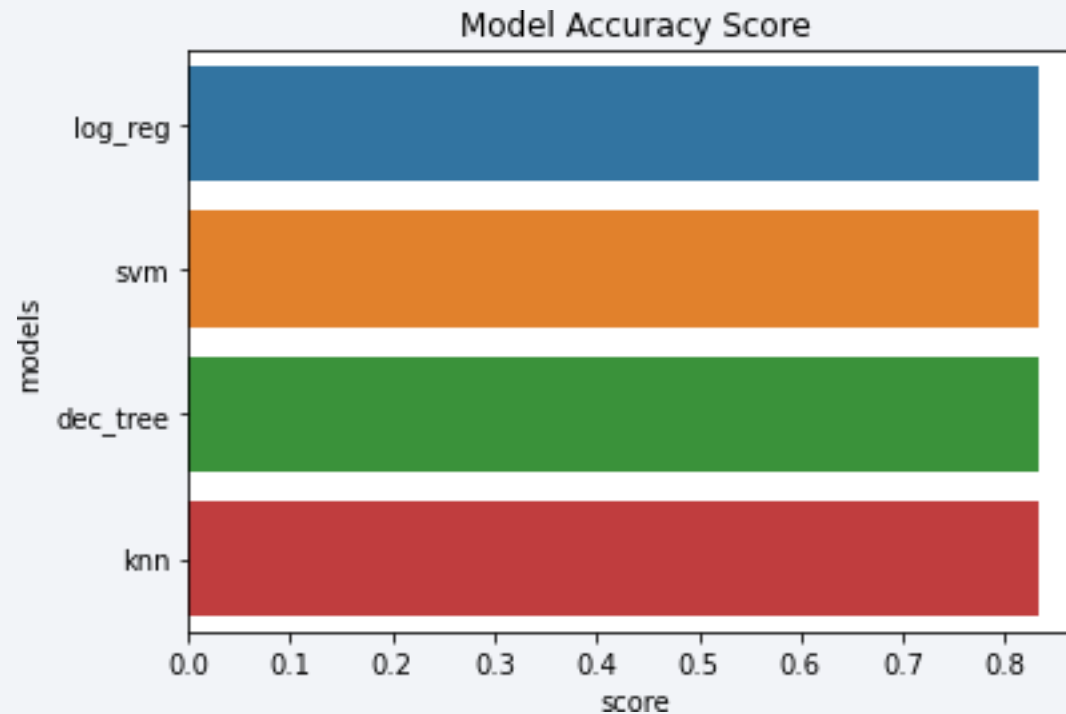
Plotly dashboard has a Payload range selector. However, this is set from 0-10000 instead of the max Payload of 15600. Class indicates 1 for successful landing and 0 for failure. Scatter plot also accounts for booster version category in color and number of launches in point size. In this particular range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

Section 5

# Predictive Analysis (Classification)

# Classification Accuracy

---



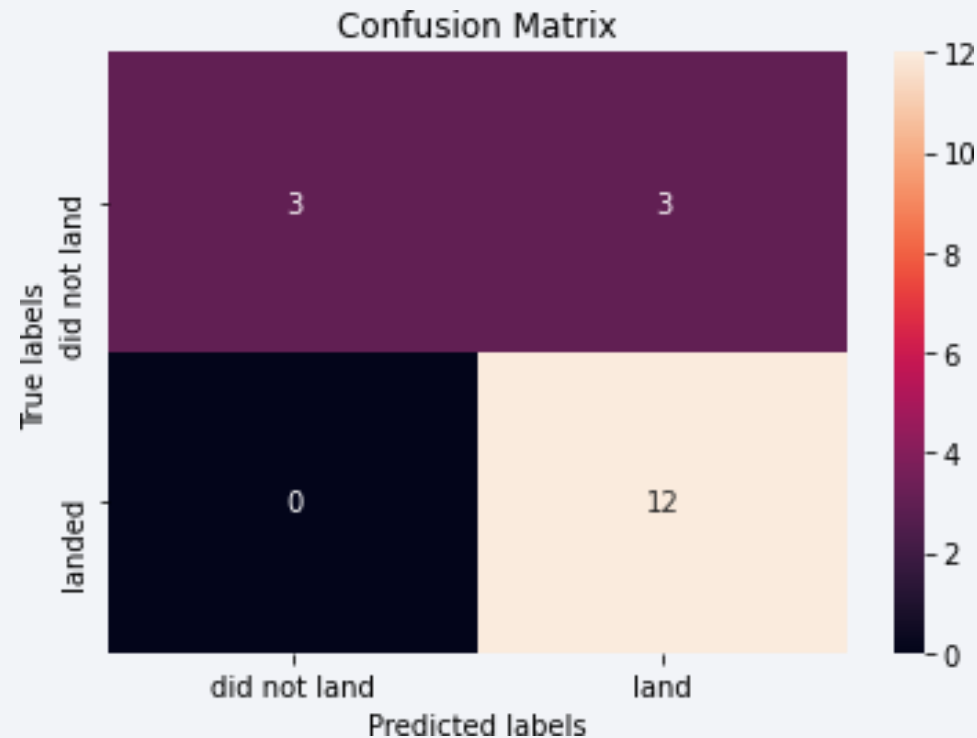
All models had virtually the same accuracy on the test set at 83.33% accuracy.

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



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.



# Conclusions

---

- 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

---

- Include any relevant assets like Python code snippets, SQL queries, charts, Notebook outputs, or data sets that you may have created during this project

Thank you!

