



Winning Space Race with Data Science

Sudheer Chowdary Pulapa
22/01/2024



Outline

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

Executive Summary

- **Objective:** Predict the success of Falcon 9 first stage landing to determine launch cost viability.
- **Data Collection and Wrangling:**
 - Utilized RESTful API and Web Scraping for Falcon 9 first-stage landing data.
 - Converted data into a Pandas data frame and performed essential data wrangling.
- **Exploratory Data Analysis (EDA):**
 - Created scatter plots and bar charts using Pandas for insightful data analysis.
 - Executed SQL queries to select and sort data, enhancing exploratory analysis.
- **Data Visualization:**
 - Developed an interactive dashboard with Plotly Dash, incorporating pie charts and scatter plots.
 - Utilized Folium to build an interactive map for analyzing launch site proximity.
- **Machine Learning for Predictive Analysis:**
 - Split data into training and testing sets for machine learning models.
 - Trained classification models (SVM, Classification Trees, Logistic Regression).
 - Conducted Hyperparameter grid search for optimization.
- **Outcome Analysis:**
 - Determined the best-performing method using test data.
 - Demonstrated the application of machine learning in predicting Falcon 9 landing success.

Introduction

- Commercial space travel is thriving, with companies like SpaceX making it more affordable.
- SpaceX stands out for its cost-effective Falcon 9 rocket launches, priced at \$62 million, a significant saving compared to other providers charging upwards of \$165 million.
- The key to SpaceX's cost efficiency lies in the reusable first stage of the Falcon 9.
- Our goal is to determine the launch price by predicting the first stage's reusability, thereby providing valuable insights for cost estimation.
- This presentation aims to gather pertinent information about SpaceX and create dashboards for our team to analyze and make informed decisions.



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 collection involved utilizing a RESTful API and web scraping methods.
- The RESTful API was employed to gather structured data, while web scraping helped extract additional relevant information.
- The collected data was converted into a data frame for efficient analysis.
- Subsequent data wrangling techniques were applied to enhance the quality and usability of the dataset.

Data Collection – SpaceX API

- The SpaceX REST API endpoints, or URL, starts with `api.spacexdata.com/v4/`
- Conducting a GET request with the request library to retrieve and utilize launch data
- Our response will be in the form of a list of JSON objects
- The link to the notebook is [Data Collection](#)

```
url="https://api.spacexdata.com/v4/launches/past"  
response =requests.get(url)  
response.json()
```


Data Collection - Scraping

- Employ the Python BeautifulSoup package for web scraping HTML containing Falcon 9 launch records.
- Followed by parsing the table data and transforming it into a Pandas data frame.
- The link to the notebook is [Data Collection Web Scraping](#)

2020 - 10/1

On 10/1/2020, SpaceX launched the first Falcon Heavy for its first of 24 launches for Starlink satellites in 2020.^[1] In addition to 14 or 15 non-Starlink launches, 40 Starlink satellites, 18 of which for Starlink satellites, Falcon 9 had the most payload ever, and 3 Falcon Heavy were launched from the same launch complex in 2020, only before SpaceX's Long March rocket family.^[2]

Flight No.	Date and Time (UTC)	Version, Revision	Launch site	Payload ^[1]	Payload mass	Orbit	Customer	Launch outcome	Recovery status
18	7 January 2020, 00:10:17 UTC	F9 90.1, 01000.0	CCAFS, SLC-40	Starlink 2 v1.2 (20 satellites)	15,800 kg (34,450 lb) ^[3]	LEO	SpaceX	Success	Success (one day)
19	19 January 2020, 10:00:00 UTC	F9 90.1, 01000.0	Vandenberg, SLC-3E	Orion EUS-1 (1 satellite)	12,500 kg (27,500 lb)	Sub-orbital	ESA	Success	No attempt
20	28 January 2020, 14:00:00 UTC	F9 90.1, 01000.0	CCAFS, SLC-40	Starlink 2 v1.2 (20 satellites)	15,800 kg (34,450 lb) ^[3]	LEO	SpaceX	Success	Success (one day)
21	17 February 2020, 14:00:00 UTC	F9 90.1, 01000.0	CCAFS, SLC-40	Starlink 2 v1.2 (20 satellites)	15,800 kg (34,450 lb) ^[3]	LEO	SpaceX	Success	Success (one day)
22	7 March 2020, 00:00:00 UTC	F9 90.1, 01000.0	CCAFS, SLC-40	Starlink 2 v1.2 (20 satellites)	15,800 kg (34,450 lb) ^[3]	LEO	SpaceX	Success	Success (one day)
23	10 March 2020, 10:00:00 UTC	F9 90.1, 01000.0	CCAFS, SLC-40	Starlink 2 v1.2 (20 satellites)	15,800 kg (34,450 lb) ^[3]	LEO	SpaceX	Success	Success (one day)

Web scraping with BeautifulSoup

Flight No.	PayloadMass	Orbit	LaunchSite	Outcome	Flights	GridFins	Reused	Legs	LandingPad	Block	ReusedCount	Serial	Index
on 1	20.0	LEO	Kwajalein Atoll	None None	1	False	False	False	None	NaN	0	Merlin1A	1
on 1	NaN	LEO	Kwajalein Atoll	None None	1	False	False	False	None	NaN	0	Merlin2A	1
on 1	165.0	LEO	Kwajalein Atoll	None None	1	False	False	False	None	NaN	0	Merlin2C	1
on 1	200.0	LEO	Kwajalein Atoll	None None	1	False	False	False	None	NaN	0	Merlin3C	1
on 9	NaN	LEO	CCAFS SLC 40	None None	1	False	False	False	None	1.0	0	B0003	-8

Data Wrangling

- We would like landing outcomes to be converted into classes y (either 1 or 0)
 - 0 is a bad outcome i.e., the booster did not land
 - 1 is a good outcome i.e., the booster did land
- Calculated the count of launches at each site and analyzed the frequency of each orbit type and its occurrences.
- The link to the notebook is [Data Wrangling](#)

EDA with Data Visualization

- Some attributes can be used to determine if the first stage can be reused.
- We want to determine what attributes are correlated with successful landings.
- Generated various plots, such as scatter points, bar charts, and line graphs, to examine the impact of one variable on another, such as FlightNumber vs. PayloadMass or LaunchSuccess's Yearly Trend.
- The link to the notebook is [EDA with Data Visualization](#)

EDA with SQL

- Using the SQL queries collectively provide insights into various aspects of space missions such as:
 - The names of unique launch sites in the space mission
 - The total payload mass carried by booster launched by NASA (CRS)
 - The average payload mass carries by booster version F9 v1.1
 - The total number of successful and failure mission outcomes
 - The failed landing outcomes in drone ship, their booster version and launch site names
- The link to the notebook is [EDA with SQL Queries](#)

Build an Interactive Map with Folium

- Utilized Folium to create an interactive map to analyze the Launch Site's Geo and Proximities empowering users to make informed decisions through interactive visual analytics.
- Added markers, circles, and lines for launch site locations and proximities.
- Provided answers by utilizing distances between a launch site and its surroundings:
 - Are launch sites near railways, highways and coastlines
 - Do launch sites keep certain distance away from cities
- The link to the notebook is [Interactive Map with Folium](#)

Build a Dashboard with Plotly Dash

- Constructed an interactive dashboard using Plotly Dash.
- Created pie charts to visualize total launches for specific launch sites.
- Generated scatter graphs to illustrate the relationship between outcome and payload mass for various boosters.
- The link to the notebook is [Plotly Dashboard](#).

Predictive Analysis (Classification)

- Loaded data using NumPy and Pandas, transformed it, and then partitioned it into training and testing sets.
- Developed various machine learning models, fine-tuning hyperparameters through GridSearchCV.
- Employed accuracy metrics for model evaluation, identifying the best-performing model. Enhanced model performance through feature engineering and algorithm tuning.
- The link to the notebook is [Predictive Analytics](#).

Results

- Exploratory data analysis results
- Interactive analytics demo in screenshots
- Predictive analysis results



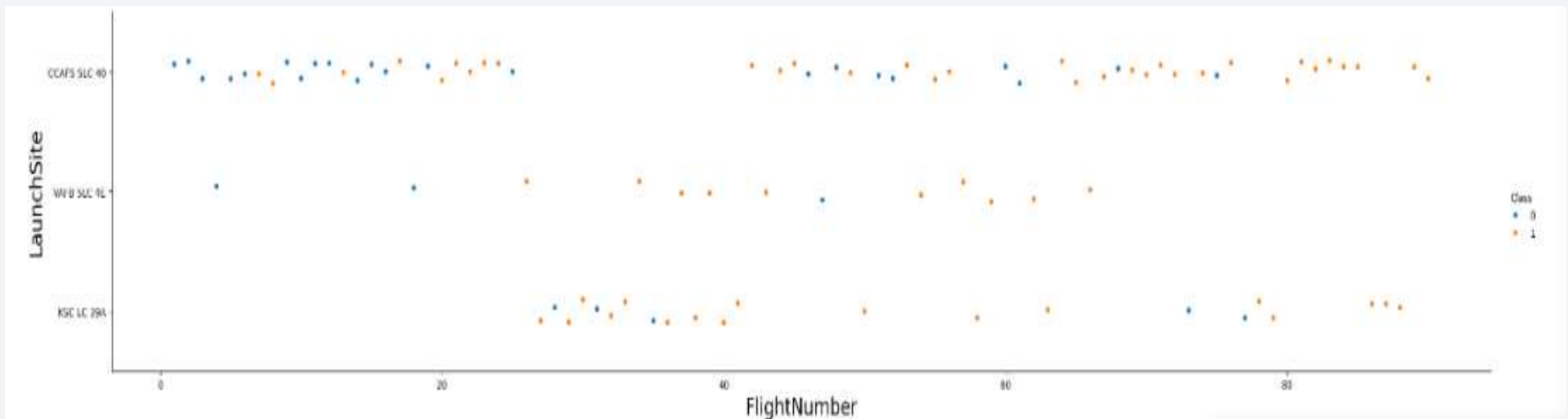
Section 2

Insights drawn from EDA

Flight Number vs. Launch Site



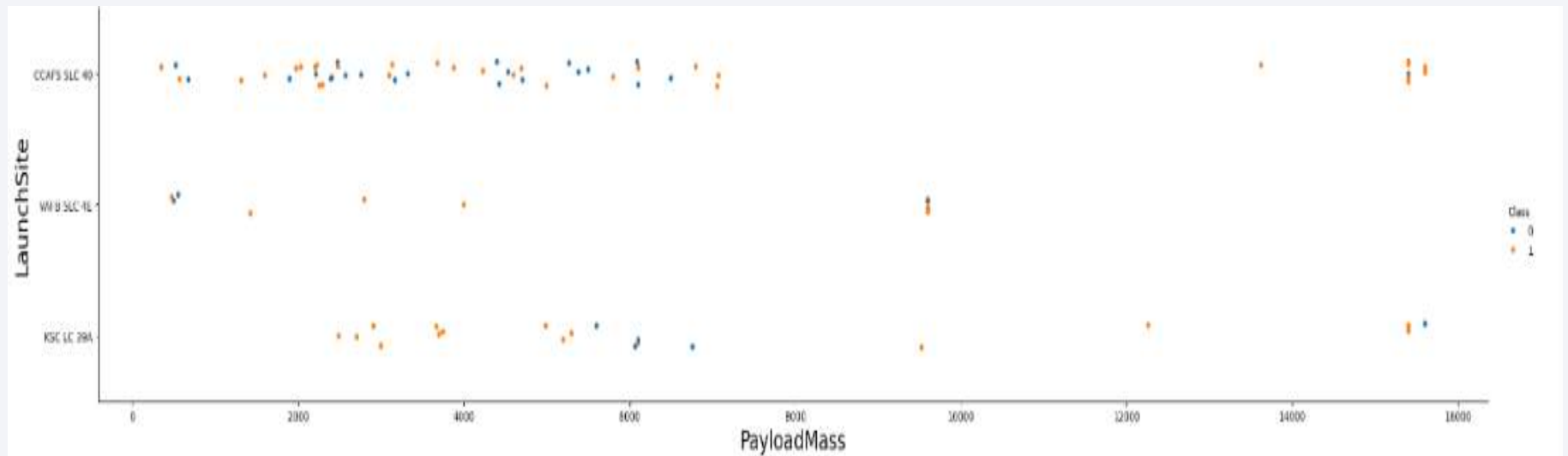
The greater the flight amount at a launch site, the greater the success rate.



Payload vs. Launch Site



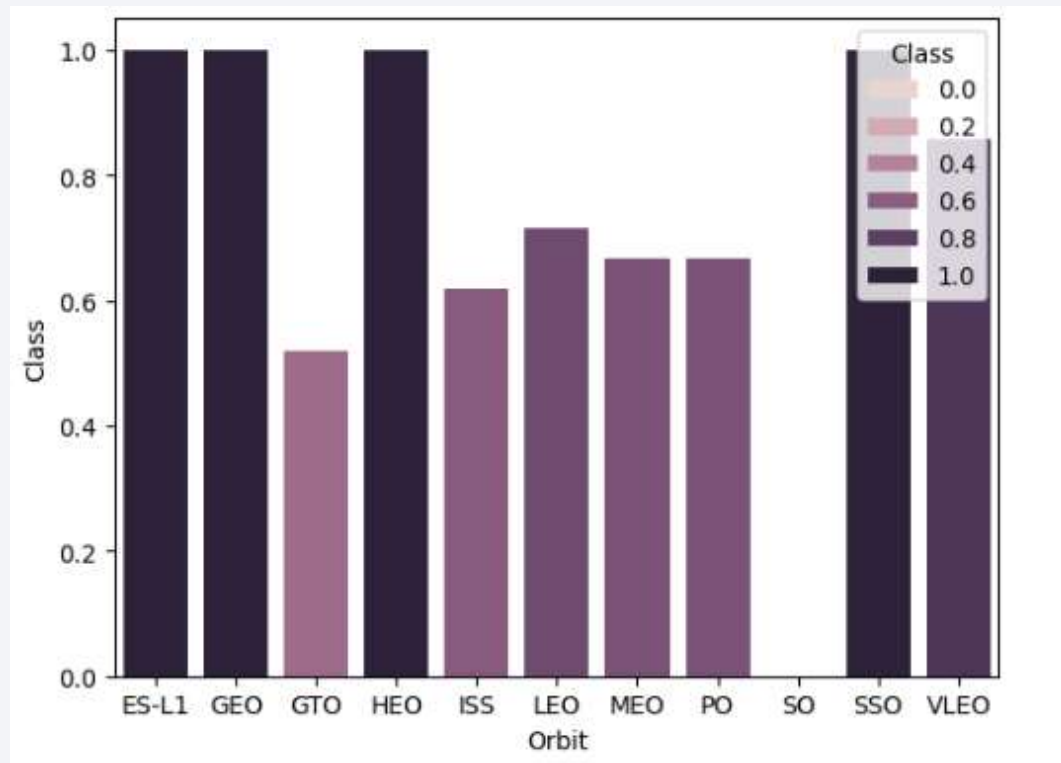
The greater the payload mass,
the higher the success rate.



Success Rate vs. Orbit Type



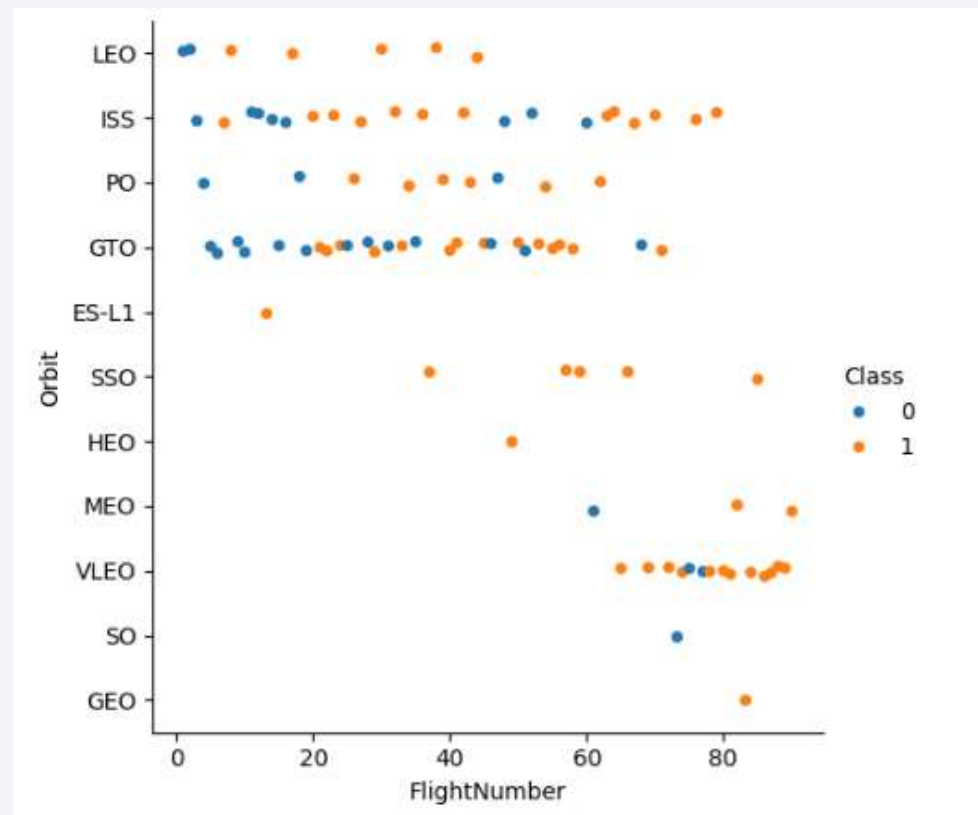
ES-L1, GEO, HEO & SSO achieved the highest success rate.



Flight Number vs. Orbit Type



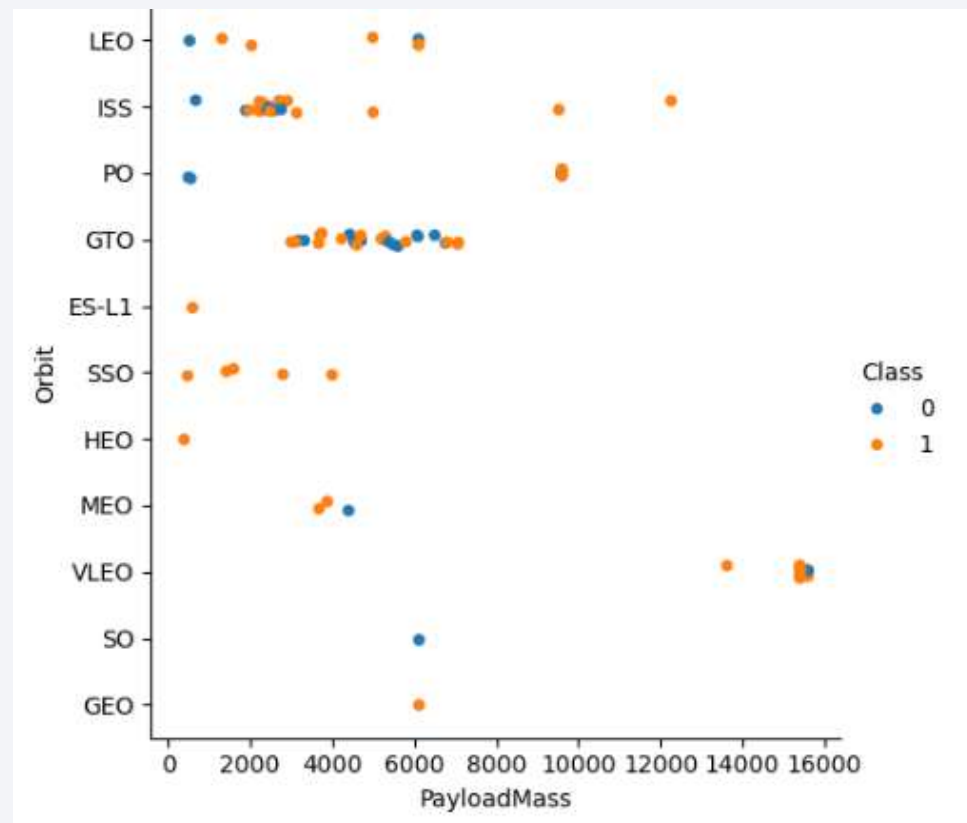
No relationship is found
between Flight Number
vs Orbit Types.



Payload vs. Orbit Type



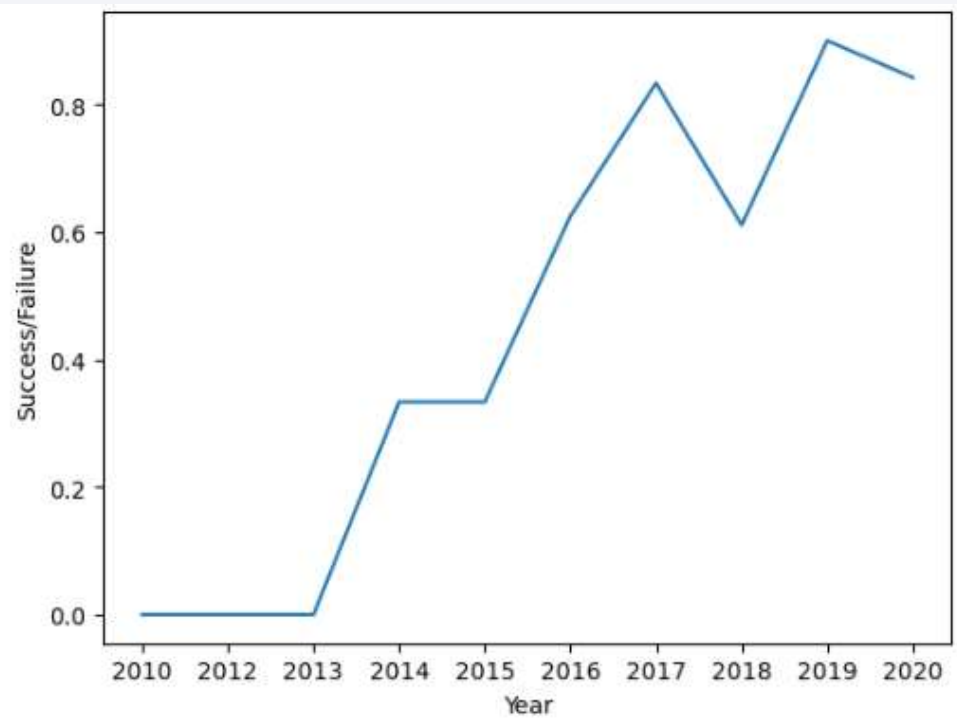
Most payloads are assigned for VLEO, PO and ISS with successful results.



Launch Success Yearly Trend



Launch site success rate increases every year since 2013.



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 SPACEXTABLE;
```

```
* 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 the keywords LIKE and LIMIT to filter launch site's beginning with 'CCA' and restricts the number of rows returned

```
%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_MASS_KG_	Orbit	Customer	Mission_Outcome	Landing_Outcome
2010-06-04	18:45:00	F9 v1.0 B0003	CCAFS LC-40	Dragon Spacecraft Qualification Unit	0	LEO	SpaceX	Success	Failure (parachute)
2010-12-08	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	7:44:00	F9 v1.0 B0005	CCAFS LC-40	Dragon demo flight C2	525	LEO (ISS)	NASA (COTS)	Success	No attempt
2012-10-08	0:35:00	F9 v1.0 B0006	CCAFS LC-40	SpaceX CRS-1	500	LEO (ISS)	NASA (CRS)	Success	No attempt
2013-03-01	15:10:00	F9 v1.0 B0007	CCAFS LC-40	SpaceX CRS-2	677	LEO (ISS)	NASA (CRS)	Success	No attempt

Total Payload Mass

Calculated the total payload mass carried by NASA's booster which resulted in 45596 kg in total.

```
%sql SELECT SUM(PAYLOAD_MASS_KG_) FROM SPACEXTABLE WHERE "Customer" LIKE "NASA (CRS)";
```

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

```
Done.
```

SUM(PAYLOAD_MASS_KG_)

45596

Average Payload Mass by F9 v1.1

Calculated the average payload mass carried by booster version F9 v1.1 which resulted in 2928.4 kg in total.

```
%sql SELECT AVG(PAYLOAD_MASS_KG_) FROM SPACEXTABLE WHERE "Booster_Version" LIKE "F9 v1.1";
```

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

```
Done.
```

AVG(PAYLOAD_MASS_KG_)

2928.4

First Successful Ground Landing Date

The first successful landing outcome on ground pad took place on 22nd December 2015.

```
%sql SELECT MIN("Date") FROM SPACE_TABLE WHERE "Landing_Outcome" LIKE "Success (ground pad)";
```

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

```
Done.
```

MIN("Date")

2015-12-22

Successful Drone Ship Landing with Payload between 4000 and 6000

Using the WHERE clause to filter successful drone ship landing AND payload mass greater than 4000 and less than 6000.

```
%sql SELECT "Booster_Version" FROM SPACEXTABLE WHERE "Landing_Outcome" LIKE "Success (drone ship)" AND "PAYLOAD_MASS__KG_" > 4000 AND "PAYLOAD_MASS__KG_" < 6000
* sqlite:///my_data1.db
Done.
```

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

Total Number of Successful and Failure Mission Outcomes

A GROUP BY clause was used to count mission outcome results.

```
%sql SELECT MISSION_OUTCOME, COUNT(*) as total_number FROM SPACEXTBL GROUP BY MISSION_OUTCOME;
* sqlite:///my_data1.db
Done.
```

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

Boosters Carried Maximum Payload

A subquery in the WHERE clause selects the maximum payload mass using the MAX function.

```
%%sql
SELECT "Booster_Version"
FROM SPACEXTABLE WHERE
PAYLOAD_MASS_KG_ = (SELECT MAX(PAYLOAD_MASS_KG_) FROM SPACEXTABLE);

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

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

2015 Launch Records

A substr() method was used to list out the details for failed launches on drone ship for the year 2015,

```
%%sql
SELECT substr("Date",6,2) as Month, "Date", "Booster_Version", "Launch_Site", "Landing_Outcome"
FROM SPACEXTABLE
WHERE "Landing_Outcome" LIKE 'Failure (drone ship)' AND substr("Date",0,5) LIKE "2015";
```

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

```
Done.
```

Month	Date	Booster_Version	Launch_Site	Landing_Outcome
01	2015-01-10	F9 v1.1 B1012	CCAFS LC-40	Failure (drone ship)
04	2015-04-14	F9 v1.1 B1015	CCAFS LC-40	Failure (drone ship)

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

We selected Landing outcomes and the COUNT of landing outcomes from the data using a WHERE clause to filter outcomes BETWEEN 2010-06-04 and 2010-03-20 with a GROUPBY clause to group each landing outcome.

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

* sqlite:///my_data1.db
Done.

Landing_Outcome	COUNT("Landing_Outcome")	Date
No attempt	10	2012-05-22
Success (drone ship)	5	2016-04-08
Failure (drone ship)	5	2015-01-10
Success (ground pad)	3	2015-12-22
Controlled (ocean)	3	2014-04-18
Uncontrolled (ocean)	2	2013-09-29
Failure (parachute)	2	2010-06-04
Precluded (drone ship)	1	2015-06-28

A satellite view of Earth from space, showing the curvature of the planet and city lights at night. The image is a composite of a solid blue rectangle on the left and a satellite photograph of Earth on the right. The Earth is shown from a high altitude, with the horizon line curving across the frame. The night side of the Earth is visible, with numerous bright yellow and orange lights from cities and towns scattered across the landmasses. The atmosphere is visible as a thin blue layer along the horizon.

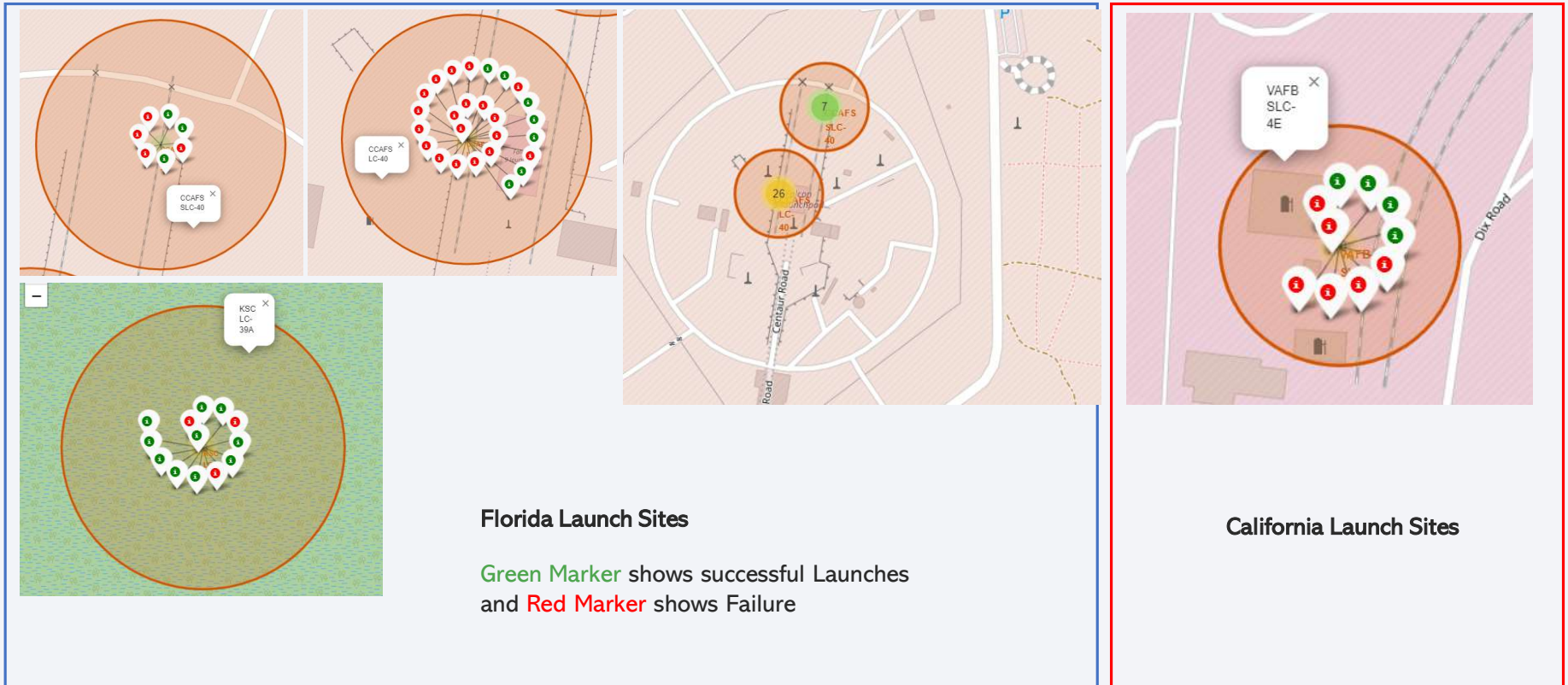
Section 3

Launch Sites Proximities Analysis

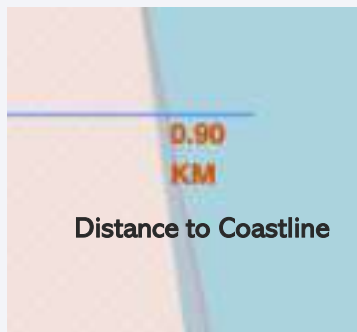
All Launch Site's Global Location



Markers Showing Launch Site Success/Failure



Launch Site Distance to Landmark



- * Are launch sites in close proximity to railways? Yes
- * Are launch sites in close proximity to highways? Yes
- * Are launch sites in close proximity to coastline? Yes
- * Do launch sites keep certain distance away from cities? No

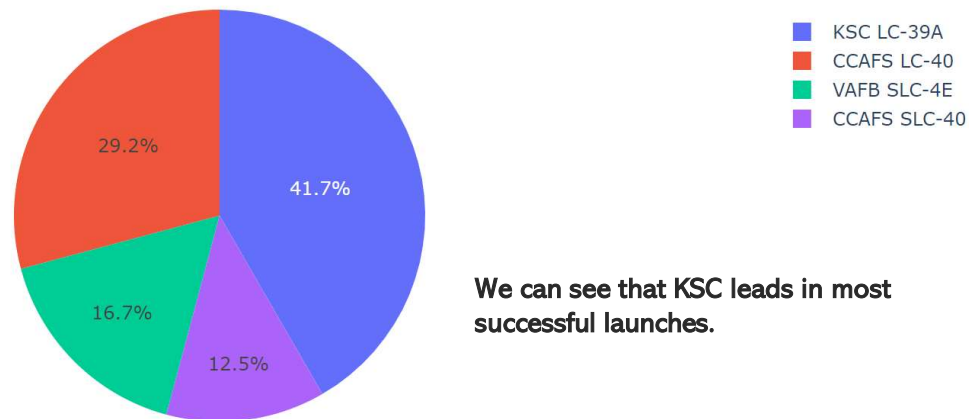


Section 4

Build a Dashboard with Plotly Dash

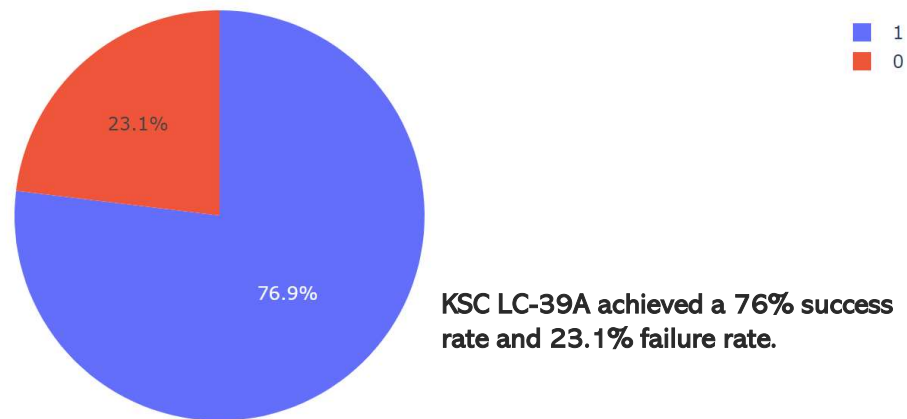
Launch Site's Success using Pie Charts

Success Count For All Launch Sites

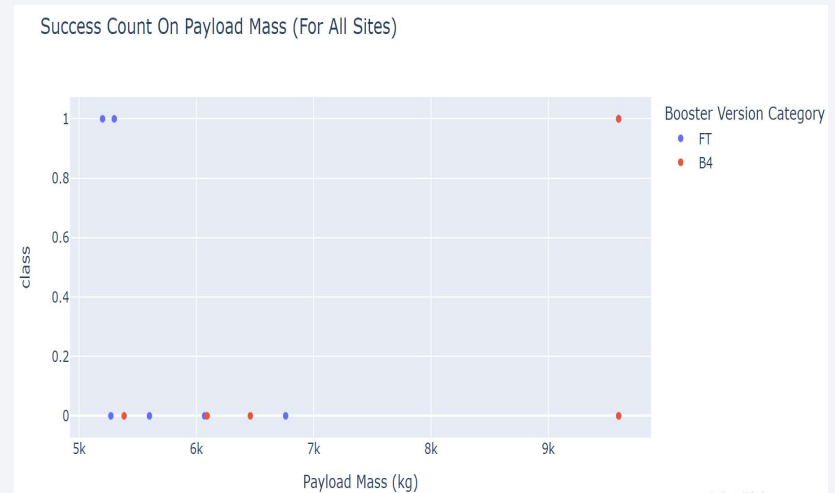
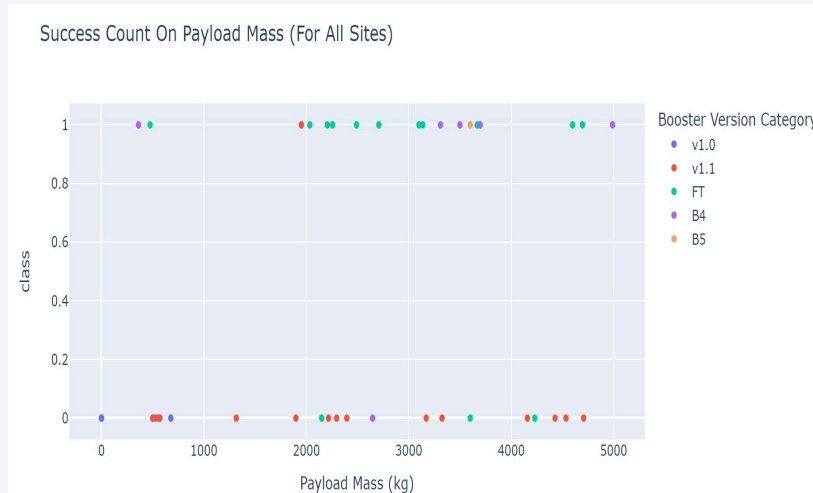


Highest Launch Success Ratio using Pie Chart

Total Success Launches For Site KSC LC-39A



Payload vs Launch Outcomes using Scatter Plot



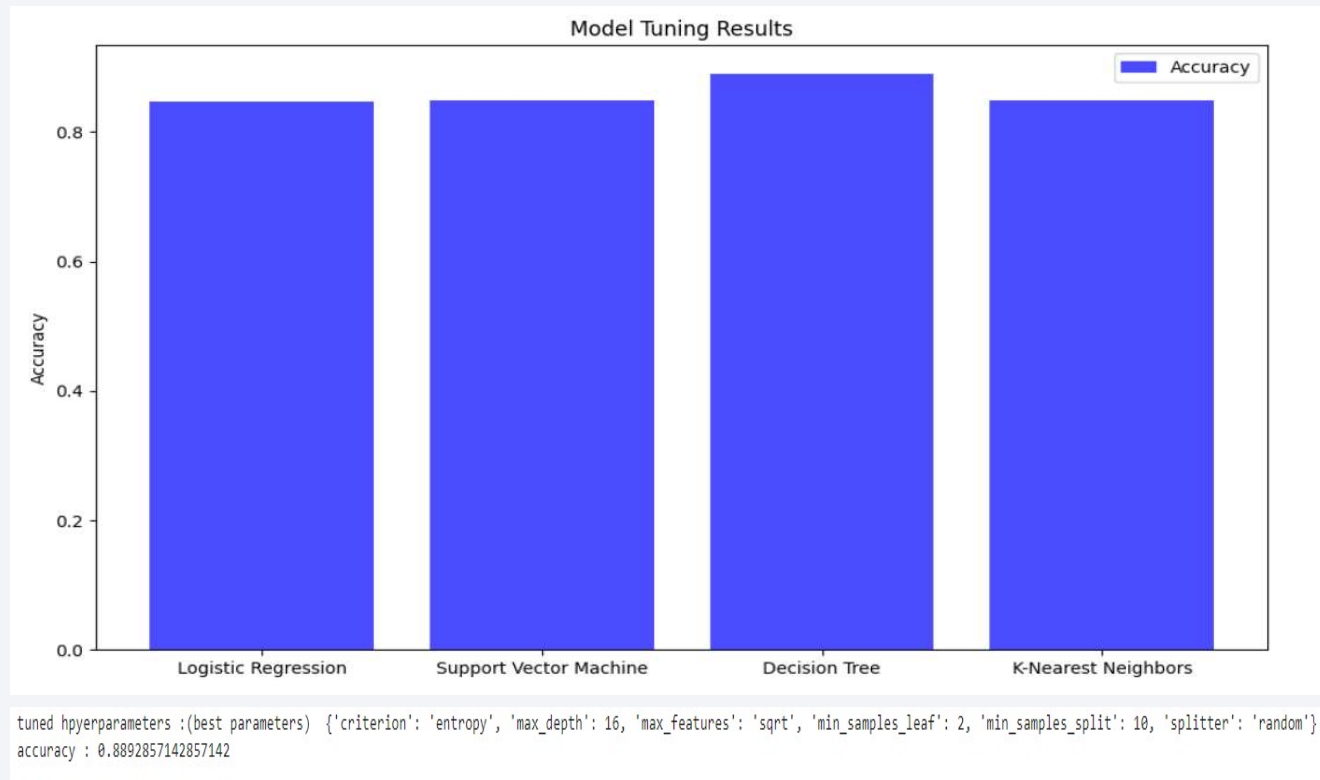
We can see the success rate for lower weighted payloads is higher than the heavier weighted payloads.



Section 5

Predictive Analysis (Classification)

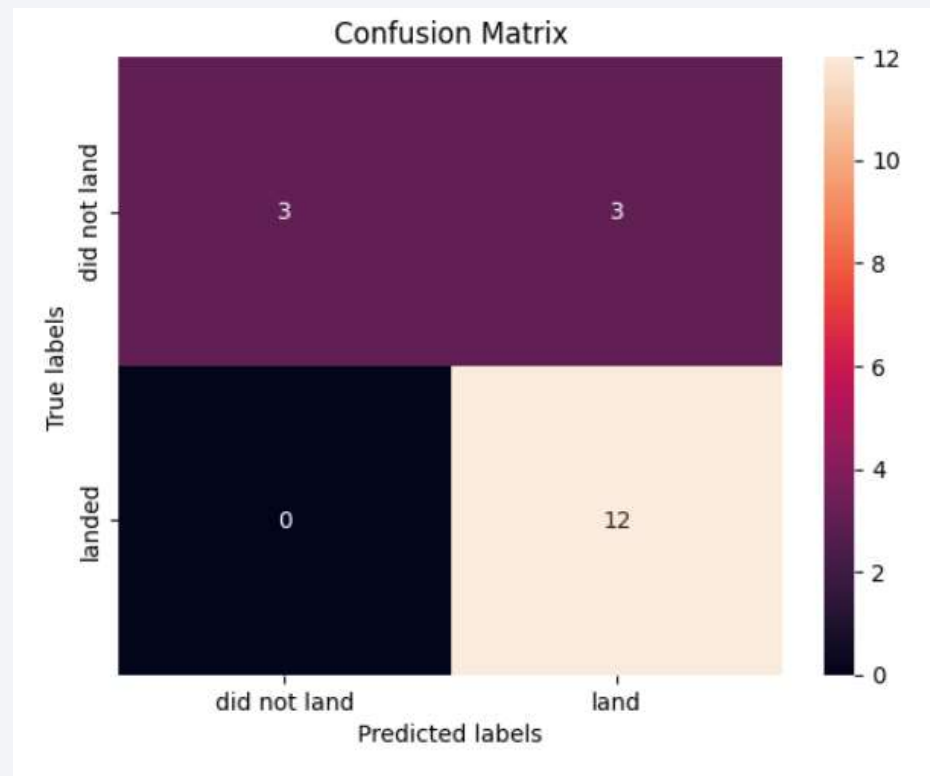
Classification Accuracy



The Decision Tree model resulted to the highest accuracy amongst the 4.

Confusion Matrix

Decision Tree Classifier Confusion Matrix reveals effective class distinction and highlights a False Positives Issue.



Conclusion

- The greater the flight amount at a launch site, the greater the success rate.
- The greater the payload mass, the higher the success rate.
- ES-L1, GEO, HEO & SSO achieved the highest success rate.
- Launch site success rate increases every year since 2013.
- The Decision Tree model resulted to the highest accuracy



Thank you!

