

Winning Space Race with Data Science

Frazer Pereira 24-08-2023

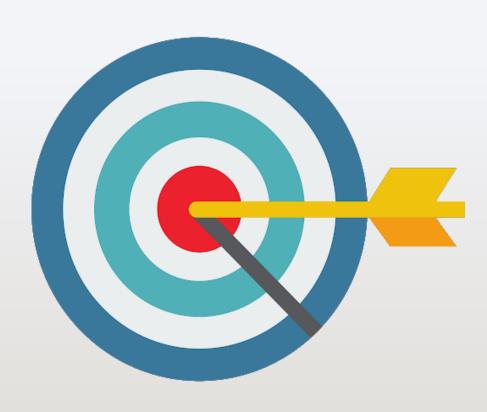


2 OUTLINE

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



Executive Summary



• Summary of methodologies



- Summary of Results
 - Exploratory Data Analysis results
 - Interactive analytics demo in screenshots
 - Predictive analysis results

4 INTRODUCTION

Project background and context:

• SpaceX is the most successful company of the commercial space age, making space travel affordable. The company 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. Based on public information and machine learning models, we are going to predict if SpaceX will reuse the first stage.

Questions to be answered:

- How do variables such as payload mass, launch site, number of flights, and orbits affect the success of the first stage landing?
- Does the rate of successful landings increase over the years?
- What is the best algorithm that can be used for binary classification in this case?





Methodology



Data Collection Methodology:

- Using SpaceX Rest API
- Using Web Scrapping from Wikipedia

Perform Data Wrangling:

- Filtering the data
- Dealing with missing values
- Using One Hot Encoding to prepare the data to a binary classification

Perform exploratory Data Analysis (EDA) using Visualization and SQL

Perform interactive visual analytics using Folium and Plotly Dash

Perform predictive analysis using classification models

Building, tuning and evaluation of classification models to ensure the best results

Data Collection process involved a combination of API requests from SpaceX REST API and Web Scraping data from a table in SpaceX's Wikipedia entry.

7 DATA COLLECTION

We had to use both of these data collection methods in order to get complete information about the launches for a more detailed analysis.

- Data Columns are obtained by using SpaceX's REST API:
- FlightNumber, Date, BoosterVersion, PayloadMass, Orbit, LaunchSite, Outcome, Flights, GridFins, Reused, Legs, LandingPad, Block, ReusedCount, Serial, Longitude, Latitude
- Data Columns are obtained by using Wikipedia Web Scraping:
- Flight No., Launch site, Payload, PayloadMass, Orbit, Customer, Launch outcome, Version Booster, Booster landing, Date, Time

Requesting rocket launch data from SpaceX API

Decoding the response content using .json()and turning it into a dataframe using .json_normalize()

Requesting needed information about the launches from SpaceX API by applying custom functions

Constructing data we have obtained into a dictionary

Exporting the data to CSV

Replacing missing values of Payload Mass column with calculated .mean() for this column

Filtering the dataframe to only include Falcon 9 launches

Creating a dataframe from the dictionary

Data Collection - Scraping

Requesting Falcon 9 launch data from Wikipedia Creating a
BeautifulSoup
object from the
HTML response

Extracting all column names from the HTML table header

Exporting the data to CSV

Creating a dataframe from the dictionary

Constructing data we have obtained into a dictionary

Collecting the data by parsing HTML tables

Data Wrangling

In the data set, there are several different cases where the booster did not land successfully. Sometimes alanding was attempted but failed due to an accident; for example, True Ocean means the mission outcome was successfully landed to a specific region of the ocean while False Ocean means the mission outcome was unsuccessfully landed to a specific region of the ocean. True RTLS means the mission outcome was successfully landed to a ground pad False RTLS means the mission outcome was unsuccessfully landed to a ground pad. True ASDS means the mission outcome was successfully landed on a drone ship False ASDS means the mission outcome was unsuccessfully landed on a drone ship.

We mainly convert those outcomes into Training Labels with "1" means the booster successfully landed, "0" means it was unsuccessful.

Perform exploratory Data Analysis and determine Training Labels

Calculate the number of launches on each site

Calculate the number and occurrence of each orbit

Calculate the number and occurrence of mission outcome per orbit type

Create alanding outcome label from Outcome column

Exporting the data to CSV

EDA with Data Visualization

Charts were plotted:

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

Scatter plots show the relationship between variables. If a relationship exists, they could be used in machine learning model.

Bar charts show comparisons among discrete categories. The goal is to show the relationship between the specific categories being compared and ameasured value.

Line charts show trends in data over time (time series).

EDA with SQL

Performed SQL queries:

- Displaying the names of the unique launch sites in the space mission
- Displaying 5 records where launch sites begin with the string 'CCA'
- Displaying the total payload mass carried by boosters launched by NASA (CRS)
- Displaying average payload mass carried by booster version F9 vl.1
- Listing the date when the first successful landing outcome in ground pad was achieved
- Listing the names of the boosters which have success in drone ship and have payload mass greater than 4000 but less than 6000
- Listing the total number of successful and failure mission outcomes
- Listing the names of the booster versions which have carried the maximum payload mass
- Listing the failed landing outcomes in drone ship, their booster versions and launch site names for the months in year 2015
- Ranking the count of landing outcomes (such as Failure (drone ship) or Success (ground pad)) between the date 2010-06-04 and 2017-03-20 in descending order

Build an Interactive Map with Folium

Markers of all Launch Sites:

- Added Marker with Circle, Popup Label and Text Label of NASA Johnson Space Center using its latitude and longitude coordinates as a start location.
- Added Markers with Circle, Popup Label and Text Label of all Launch Sites using their latitude and longitude coordinates to show their geographical locations and proximity to Equator and coasts.

Colored Markers of the launch outcomes for each Launch Site:

- Added colored Markers of success (Green) and failed (Red) launches using Marker Cluster to identify which launch sites have relatively high success rates.

Distances between a Launch Site to its proximities:

- Added coloured Lines to show distances between the Launch Site KSC LC-39A (as an example) and its proximities like Railway, Highway, Coastline and Closest City.

Build a Dashboard with Plotly Dash

Launch Sites Dropdown List:

- Added a dropdown list to enable Launch Site selection.

Pie Chart showing Success Launches (All Sites/Certain Site):

- Added a pie chart to show the total successful launches count for all sites and the Success vs. Failed counts for the site, if a specific Launch Site was selected.

Slider of Payload Mass Range:

- Added a slider to select Payload range.

Scatter Chart of Payload Mass vs. Success Rate for the different Booster Versions:

Added a scatter chart to show the correlation between Payload and Launch Success

Predictive Analysis (Classification)

Creating a NumPy array from the column "Class" in data

Standardizing the data with StandardScaler, then fitting and transforming it

Splitting the data into training and testing sets with train_test_split function

Creating a
GridSearchCV
object with cv =
10 to find the best
parameters

Finding the method performs best by examining the Jaccard_score and F1_score metrics

Examining the confusion matrix for all models

Calculating the accuracy on the test data using the method .score() for all models

Applying
GridSearchCV
on LogReg,
SVM,
Decision Tree, and
KNN models

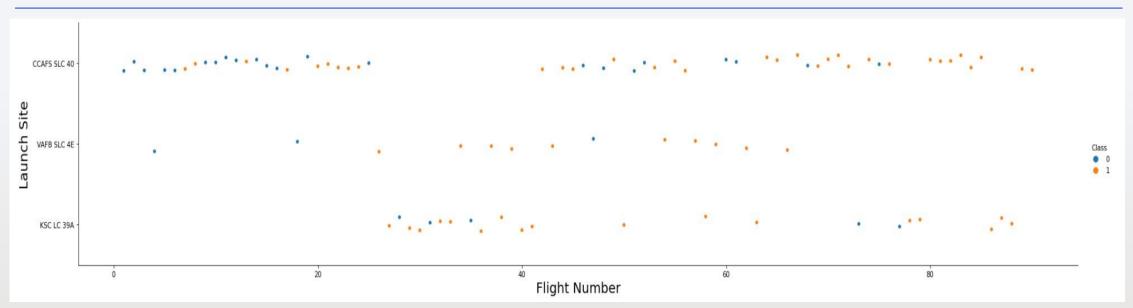
Results



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

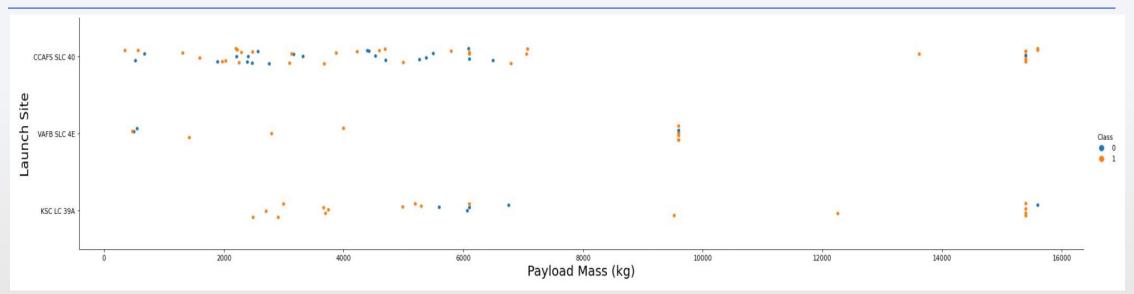


Flight Number vs. Launch Site



- The earliest flights all failed while the latest flights all succeeded.
- The CCAFS SLC 40 launch site has about a half of all launches.
- VAFB SLC 4E and KSC LC 39A have higher success rates.
- It can be assumed that each new launch has a higher rate of success.

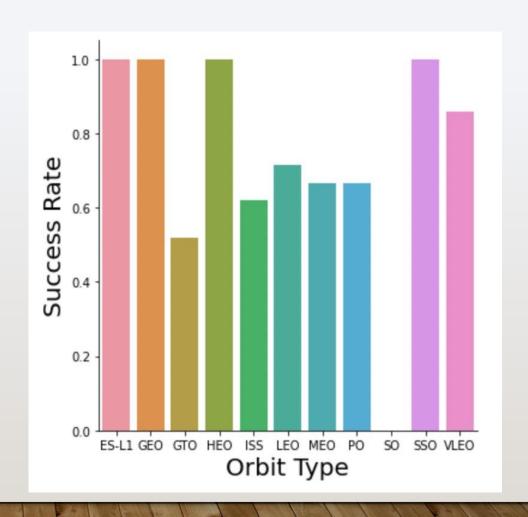
Payload vs. Launch Site



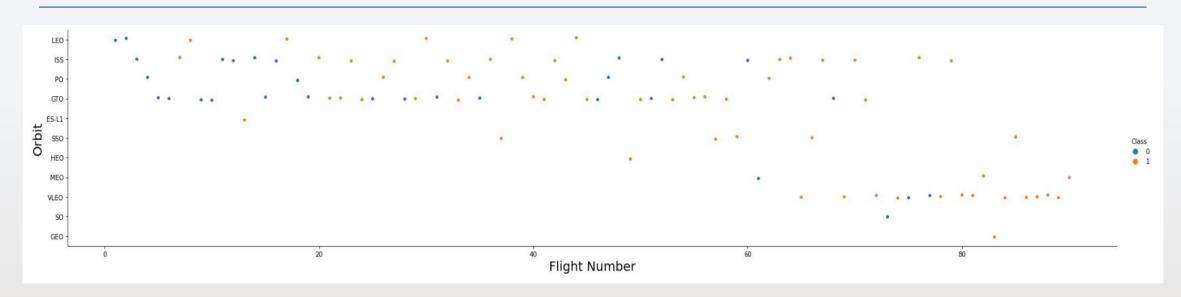
- For every launch site the higher the payload mass, the higher the success rate.
- Most of the launches with payload mass over 7000 kg were successful.

Success Rate vs. Orbit Type

- Orbits with 100% success rate:
 - ES-L1, GEO, HEO, SSO
- Orbits with 0% success rate:
 - SO
- Orbits with success rate between 50% and 85%:
 - GTO, ISS, LEO, MEO, PO



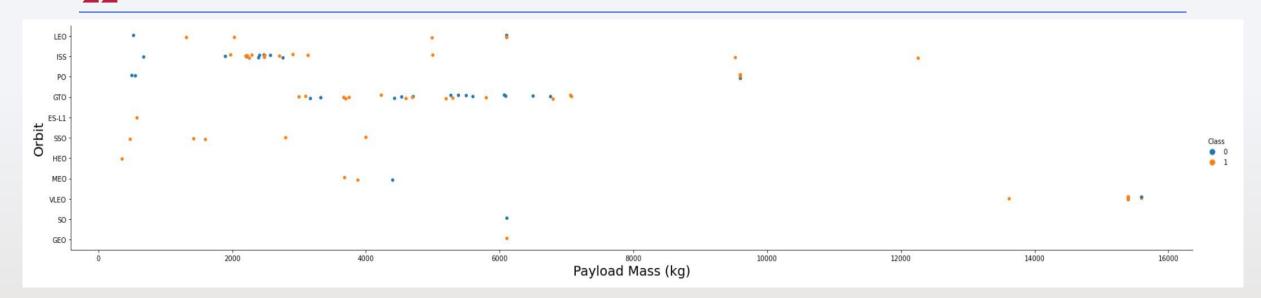
21



Explanation:

• In the LEO orbit the Success appears related to the number of flights; on the other hand, there seems to be no relationship between flight number when in GTO orbit.

Payload vs. Orbit Type



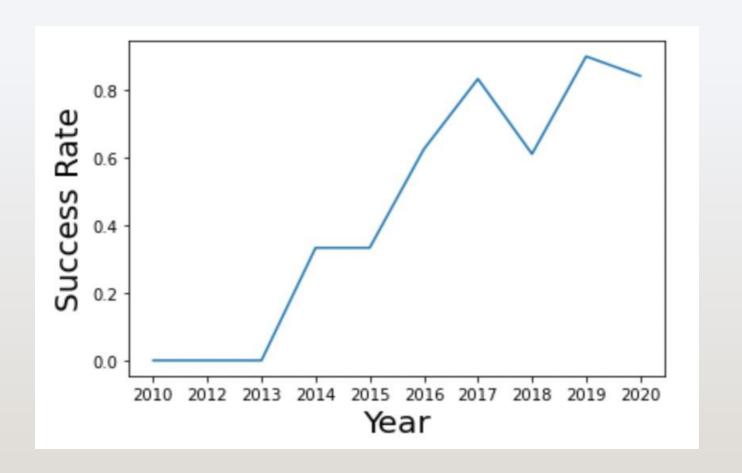
Explanation:

• Heavy payloads have a negative influence on GTO orbits and positive on GTO and Polar LEO (ISS) orbits.

Launch Success Yearly Trend

Explanation:

• The success rate since 2013 kept increasing till 2020.



All Launch Site Names

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

[11]: sql SELECT DISTINCT LAUNCH_SITE FROM SPACEXTBL ORDER BY 1:

* sqlite:///my_datal.db
Done.

[11]: Launch_Site

CCAFS LC-40

CCAFS SLC-40

KSC LC-39A

VAFB SLC-4E
```

Explanation:

• Displaying the names of the unique launch sites in the space mission.

Launch Site Names Begin with 'CCA'

12]:	sql SE	LECT * FI	ROM SPACEXTBL WH	IERE LAUNCH_S	SITE LIKE 'CO	CA%' LIMIT 5;		· 1	↓
	* sqlite:///my_data1.db Done.								
12]:	Date	Time (UTC)	Booster_Version	Launch_Site	Payload	PAYLOAD_MASSKG_	Orbit	Customer	Mission_Outcome
	2010- 04-06	18:45:00	F9 v1.0 B0003	CCAFS LC- 40	Dragon Spacecraft Qualification Unit	0	LEO	SpaceX	Success
	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
	2012- 05-22	07:44:00	F9 v1.0 B0005	CCAFS LC- 40	Dragon demo flight C2	525	LEO (ISS)	NASA (COTS)	Success
	2012- 08-10	00:35:00	F9 v1.0 B0006	CCAFS LC- 40	SpaceX CRS-1	500	LEO (ISS)	NASA (CRS)	Success

Explanation:

• Displaying records where launch sites begin with the string 'CCA'.

Total Payload Mass

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

sql SELECT SUM(PAYLOAD_MASS__KG_) AS TOTAL_PAYLOAD FROM SPACEXTBL WHERE PAYLOAD LIKE '%CRS%';

* sqlite:///my_datal.db
Done.

TOTAL_PAYLOAD

111268
```

Explanation:

• Displaying the total payload mass carried by boosters launched by NASA (CRS).

Average Payload Mass by F9 v1.1

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

sql SELECT AVG(PAYLOAD_MASS__KG_) AS AVG_PAYLOAD FROM SPACEXTBL WHERE BOOSTER_VERSION = 'F9 v1.1';

* sqlite://my_datal.db
Done.

AVG_PAYLOAD

2928.4
```

Explanation:

• Displaying average payload mass carried by booster version F9 vl.l.

First Successful Ground Landing Date

```
sql SELECT MIN(DATE) AS FIRST_SUCCESS_GP FROM SPACEXTBL WHERE LANDING_OUTCOME =_'Success_(ground pad)';

* sqlite://my_data1.db
Done.

FIRST_SUCCESS_GP

2015-12-22
```

Explanation:

• Listing the date when the first successful landing outcome in ground pad was achieved.

Successful Drone Ship Landing with Payload between 4000 and 6000

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

sql SELECT DISTINCT BOOSTER_VERSION FROM SPACEXTBL WHERE PAYLOAD_MASS__KG__BETWEEN_4000_AND_6000_AND_LANDIN * sqlite://my_datal.db
Done.

Booster_Version

F9 FT B1022

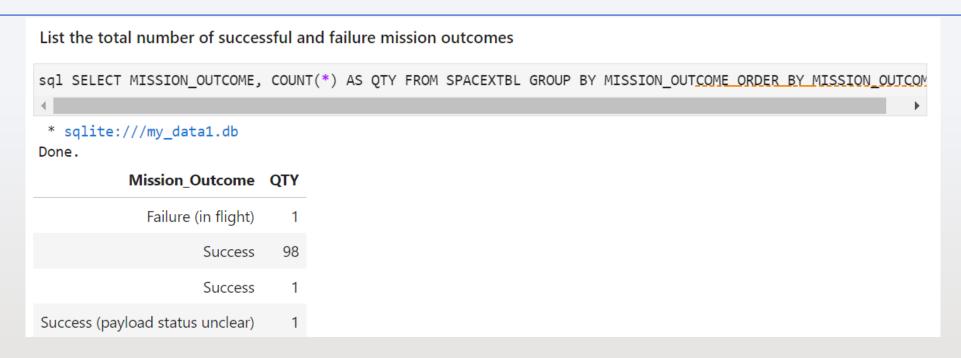
F9 FT B1021.2

F9 FT B1031.2

Explanation:

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

Total Number of Successful and Failure Mission Outcomes



Explanation:

• Listing the total number of successful and failure mission outcomes.

Boosters Carried Maximum Payload

```
List the names of the booster_versions which have carried the maximum payload mass. Use a subquery

sql SELECT DISTINCT BOOSTER_VERSION FROM SPACEXTBL WHERE PAYLOAD_MASS__KG_ = (SELECT_MAX(PAYLOAD_MASS__KG_))

* sqlite://my_data1.db
Done.

Booster_Version

F9 B5 B1048.4

F9 B5 B1049.4

F9 B5 B1049.5
```

Explanation:

• Listing the names of the booster versions which have carried the maximum payload mass.

2015 Launch Records

```
sql SELECT booster_version, launch_site, landing_outcome, date FROM spacextbl WHERE landing_outcome LIKE '%

* sqlite:///my_data1.db
Done.

Booster_Version Launch_Site Landing_Outcome Date

F9 v1.1 B1012 CCAFS LC-40 Failure (drone ship) 2015-10-01

F9 v1.1 B1015 CCAFS LC-40 Failure (drone ship) 2015-04-14

F9 v1.1 B1018 CCAFS LC-40 Precluded (drone ship) 2015-06-28
```

Explanation:

• Listing the failed landing outcomes in drone ship, their booster versions and launch site names for the months in year 2015.

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

```
sql SELECT LANDING_OUTCOME, COUNT(*) AS QTY FROM SPACEXTBL WHERE DATE BETWEEN '2010-06-04' AND '2017-03-20'

* sqlite:///my_datal.db
Done.

Landing_Outcome QTY

No attempt 10

Success (ground pad) 5

Success (drone ship) 5

Failure (drone ship) 5

Controlled (ocean) 3
```

Explanation:

• Ranking the count of landing outcomes (such as Failure (drone ship) or Success (ground pad)) between the date 2010-06-04 and 2017-03-20 in descending order.



All launch sites' location markers on a global map

Explanation:

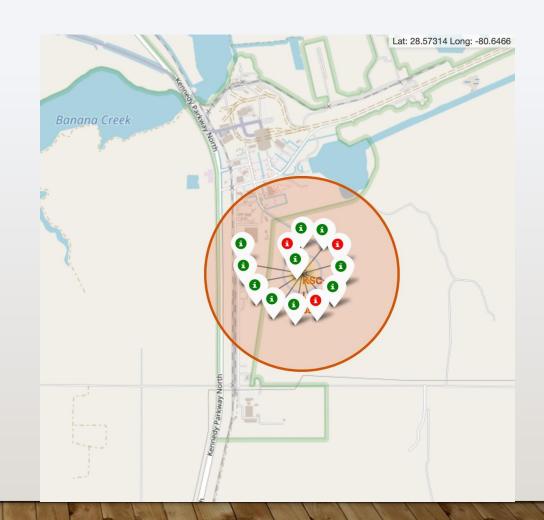
- Most of Launch sites are in proximity to the Equator line. The land is moving faster at the equator than any other place on the surface of the Earth. Anything on the surface of the Earth at the equator is already moving at 1670 km/hour. If a ship is launched from the equator it goes up into space, and it is also moving around the Earth at the same speed it was moving before launching. This is because of inertia. This speed will help the spacecraft keep up a good enough speed to stay in orbit.
- All launch sites are in very close proximity to the coast, while launching rockets towards the ocean it minimizes the risk of having any debris dropping or exploding

United States

near people:

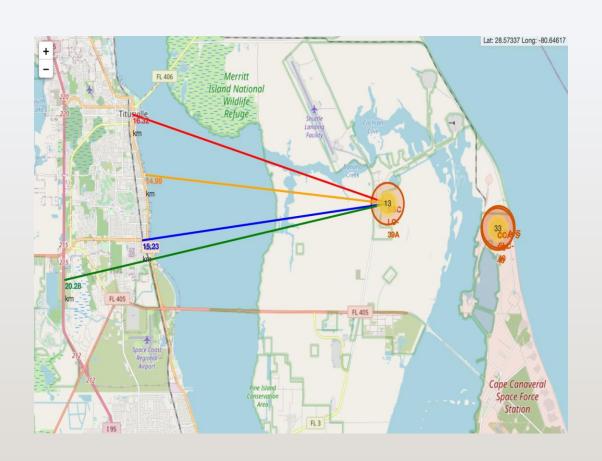
Color-labeled launch records on the map

- From the color-labeled markers we should be able to easily identify which launch sites have relatively high success rates.
 - Green Marker = Successful Launch
 - Red Marker = Failed Launch
- Launch Site KSC LC-39A has a very high Success Rate.



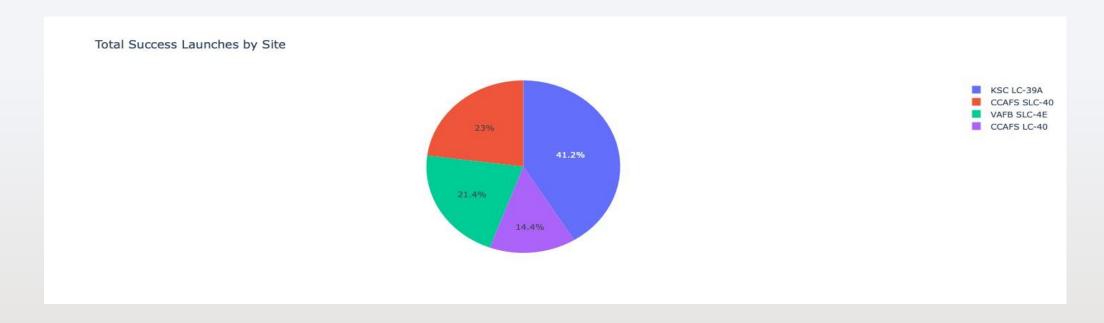
Distance from the launch site KSC LC-39A to its proximities

- From the visual analysis of the launch site KSC LC-39A we can clearly see that it is:
 - relative close to railway (15.23 km)
 - relative close to highway (20.28 km)
 - relative close to coastline (14.99 km)
- Also the launch site KSC LC-39A is relative close to its closest city Titusville (16.32 km).
- Failed rocket with its high speed can cover distances like 15-20 km in few seconds. It could be potentially





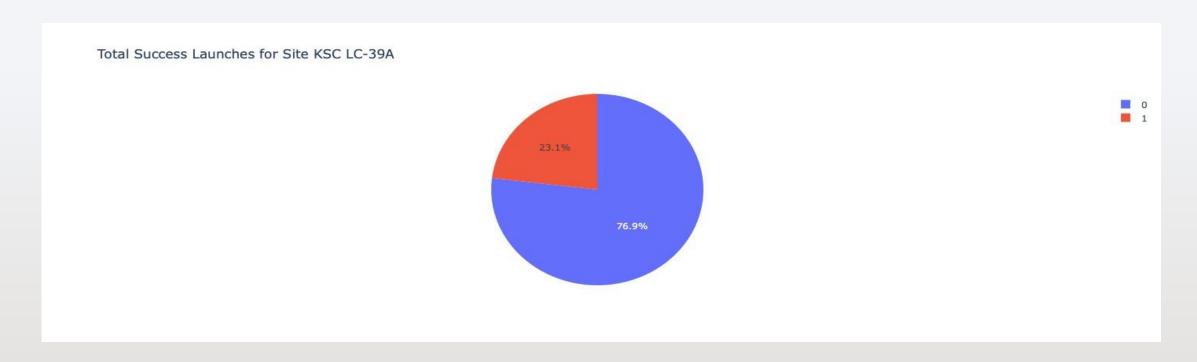
Launch success count for all sites



Explanation:

• The chart clearly shows that from all the sites, KSC LC-39A has the most successful launches.

Launch site with highest success ratio



Explanation:

• KSC LC-39A has the highest launch success rate (76.9%) with 10 successful and only 3 failed landings.

Payload Mass vs. Launch Outcome for all sites

Explanation:

• The charts show that payloads between 2000 and 5500 kg have the highest success rate.





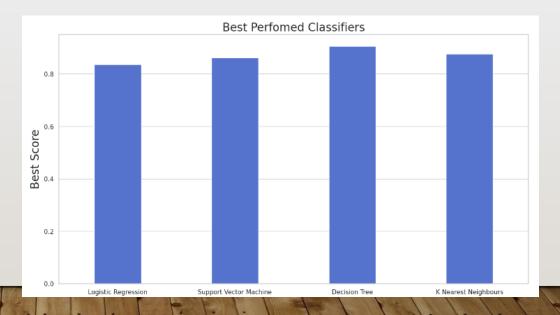


Classification Accuracy

Explanation:

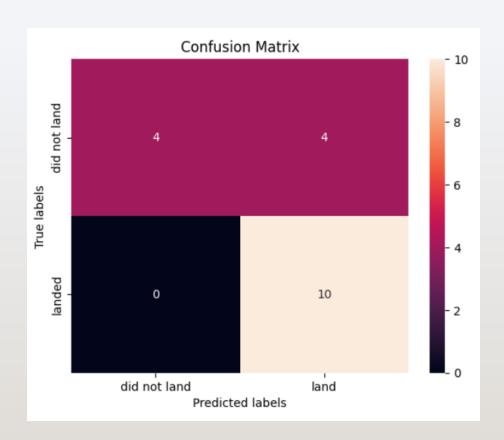
• Based on the plot and best scores, we can conclude that Decision Tree Classifier has the highest accuracy.

	Classifiers	Accuracy Score	Best Score
0	Logistic Regression	0.722222	0.835714
1	Support Vector Machine	0.777778	0.862500
2	Decision Tree	0.777778	0.905357
3	K Nearest Neighbours	0.777778	0.876786



Confusion Matrix

- The model correctly predicted 4 instances as negative and they were actually negative (True Negatives).
- The model incorrectly predicted 4 instances as positive when they were actually negative (False Positives).
- The model correctly predicted 10 instances as positive and they were actually positive (True Positives).
- The model didn't make any false negative predictions, as all instances that were actually positive were predicted as positive (False Negatives = 0).



Conclusions



- Decision Tree Model is the best algorithm for this dataset.
- Launches with a low payload mass show better results than launches with a larger payload mass.
- Most of launch sites are in proximity to the Equator line and all the sites are in very close proximity to the coast.
- The success rate of launches increases over the years.
- KSCLC-39A has the highest success rate of the launches from all the sites.
- Orbits ES-L1, GEO, HEO and SSO have 100% success rate.

