




DATA SCIENCE CAPSTONE PROJECT

Monalisa Pal

29/04/2024



OUTLINE

- Executive Summary
 - Introduction
 - Methodology
 - Data Collection and Data Wrangling
 - EDA and interactive visual analytics
 - Predictive analysis
 - EDA with visualization results
 - EDA with SQL results
 - Interactive map with Folium
 - Plotly Dash dashboard
 - Predictive analysis (classification)
 - Conclusion
- 

Executive Summary

We built a dataset to predict SpaceX launch outcomes

- We gathered information from two sources: the public SpaceX API and their Wikipedia page.
- We added a label indicating successful landings ("class" column).
- We explored the data using various techniques: SQL queries, data visualizations, interactive maps (folium), and dashboards.
- We identified the most relevant data points to use for prediction (features).
- We converted categorical data into a binary format using one-hot encoding.
- To prepare the data for machine learning, we standardized the values.
- We employed GridSearchCV to find the optimal settings for different machine learning models.
- Finally, we created visualizations to compare the accuracy of each model.

We evaluated four machine learning models (Logistic Regression, Support Vector Machine, Decision Tree Classifier, and K Nearest Neighbors) to predict successful SpaceX landings. All models achieved a similar accuracy of around 83.33%, but they tended to overestimate successful landings. To improve the models' performance and accuracy, we likely need more data.




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





Methodology

Overview of Data collection, Data Wrangling, Visualization,
Dashboard and model methods




Data Collection

- 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 Wrangling

- We combined information from two existing columns: "Mission Outcome" and "Landing Location."
 - The new "class" label classifies landings as successful (1) or unsuccessful (0) based on specific criteria:
Successful landings (1): "True ASDS," "True RTLS," and "True Ocean" in the "Mission Outcome" column.
 - Unsuccessful landings (0): Any combination of "None None," "False ASDS," "None ASDS," "False Ocean," or "False RTLS" in the "Mission Outcome" and "Landing Location" columns.
- 

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

Results:

- 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 mass appears to fall mostly between 0-6000 kg. Different launch sites also seem to use different payload mass.
- 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
- Launch Orbit preferences changed over Flight Number. Launch Outcome seems to correlate
- Success generally increases over time since 2013 with a slight dip in 2018
- Success in recent years at around 80%

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

Results:

In [4]: %%sql
SELECT UNIQUE LAUNCH_SITE
FROM SPACEXDATASET;

* ibm_db_sa://ftb12020:***@0c77d6f:
Done.

Out[4]:

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

All Launch Site Names

In [5]:


```
%%sql
SELECT *
FROM SPACEXDATASET
WHERE LAUNCH_SITE LIKE 'CCA%'
LIMIT 5;
```

* ibm_db_sa://ftb12020:***@0c77d6f2-5da9-48a9-81f8-86b520b87518.bs2io90l08kqb1od8l1cg.databases.appdomain.cloud:31198/bludb
Done.

Out[5]:

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	07:44:00	F9 v1.0 B0005	CCAFS LC-40	Dragon demo flight C2	525	LEO (ISS)	NASA (COTS)	Success	No attempt
2012-10-08	00:35:00	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

Launch Sites Name Beginning With 'CCA'




```
%%sql
SELECT SUM(PAYLOAD_MASS__KG_) AS SUM_PAYLOAD_MASS_KG
FROM SPACEXDATASET
WHERE CUSTOMER = 'NASA (CRS)';
```

```
* ibm_db_sa://ftb12020:***@0c77d6f2-5da9-48a9-81f8-86
Done.
```

sum_payload_mass_kg
45596

Total payload mass from NASA




```
%%sql
SELECT AVG(PAYLOAD_MASS_KG_) AS AVG_PAYLOAD_MASS_KG
FROM SPACEXDATASET
WHERE booster_version = 'F9 v1.1'
```

```
* ibm_db_sa://ftb12020:***@0c77d6f2-5da9-48a9-81f8-86
Done.
```

avg_payload_mass_kg
2928

Average payload mass by F9v1.1



```
%%sql
SELECT MIN(DATE) AS FIRST_SUCCESS
FROM SPACEXDATASET
WHERE landing__outcome = 'Success (ground pad)';

* ibm_db_sa://ftb12020:***@0c77d6f2-5da9-48a9-81
Done.
```

first_success

2015-12-22

First Successful Ground Pad landing Date




```
%%sql
SELECT booster_version
FROM SPACEXDATASET
WHERE landing__outcome = 'Success (drone ship)' AND payload_mass__kg_ BETWEEN 4001 AND 5999;

* ibm_db_sa://ftb12020:***@0c77d6f2-5da9-48a9-81f8-86b520b87518.bs2io90l08kqb1od8lcg.database
Done.
```

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

Successful Drone Ship Landing with Payload
Between 4000 and 6000




```
%%sql
SELECT mission_outcome, COUNT(*) AS no_outcome
FROM SPACEXDATASET
GROUP BY mission_outcome;
```

```
* ibm_db_sa://ftb12020:***@0c77d6f2-5da9-48a9-!
Done.
```

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

Total Number of Each Mission Outcome




```
%%sql
SELECT booster_version, PAYLOAD_MASS_KG_
FROM SPACEXDATASET
WHERE PAYLOAD_MASS_KG_ = (SELECT MAX(PAYLOAD_MASS_KG_) FROM SPACEXDATASET);

* ibm_db_sa://ftb12020:***@0c77d6f2-5da9-48a9-81f8-86b520b87518.bs2io90l08kqb1
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

Boosters that Carried Maximum Payload




```
%%sql
SELECT MONTHNAME(DATE) AS MONTH, landing__outcome, booster_version, PAYLOAD_MASS_KG_, launch_site
FROM SPACEXDATASET
WHERE landing__outcome = 'Failure (drone ship)' AND YEAR(DATE) = 2015;

* ibm_db_sa://ftb12020:***@0c77d6f2-5da9-48a9-81f8-86b520b87518.bs2io90l08kqb1od8lcg.databases.app
Done.
```

MONTH	landing__outcome	booster_version	payload_mass_kg_	launch_site
January	Failure (drone ship)	F9 v1.1 B1012	2395	CCAFS LC-40
April	Failure (drone ship)	F9 v1.1 B1015	1898	CCAFS LC-40

2015 Failed Drone Ship Landing Records



```
%%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;
```

```
* ibm_db_sa://ftb12020:***@0c77d6f2-5da9-48a9-81f8-86b520b87518.bs2io90l08kqb1od8lcg
Done.
```

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

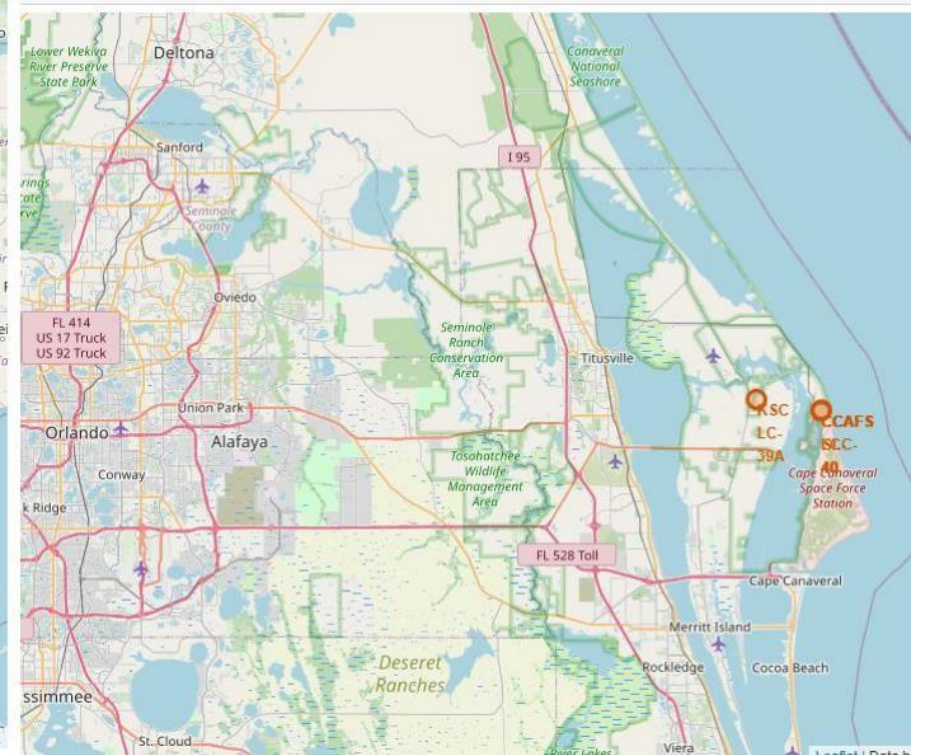
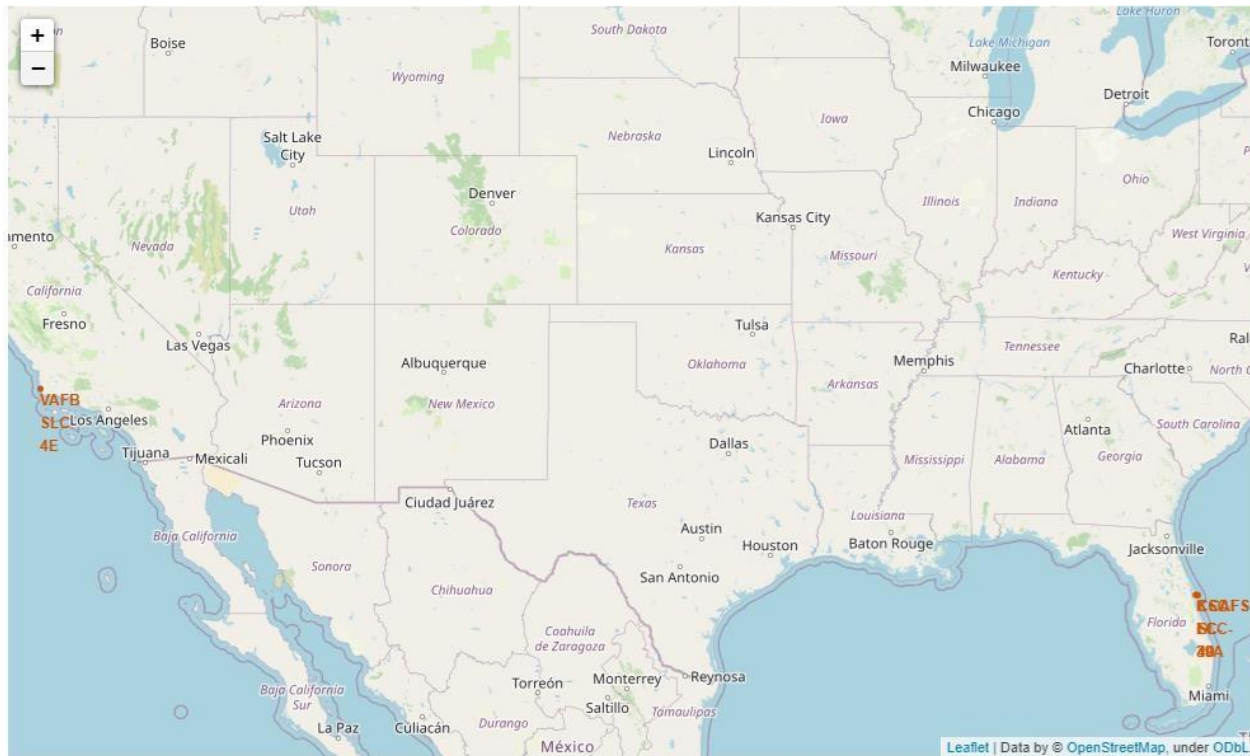
Ranking Counts of Successful Landings
Between 2010-06-04 and 2017-03-20



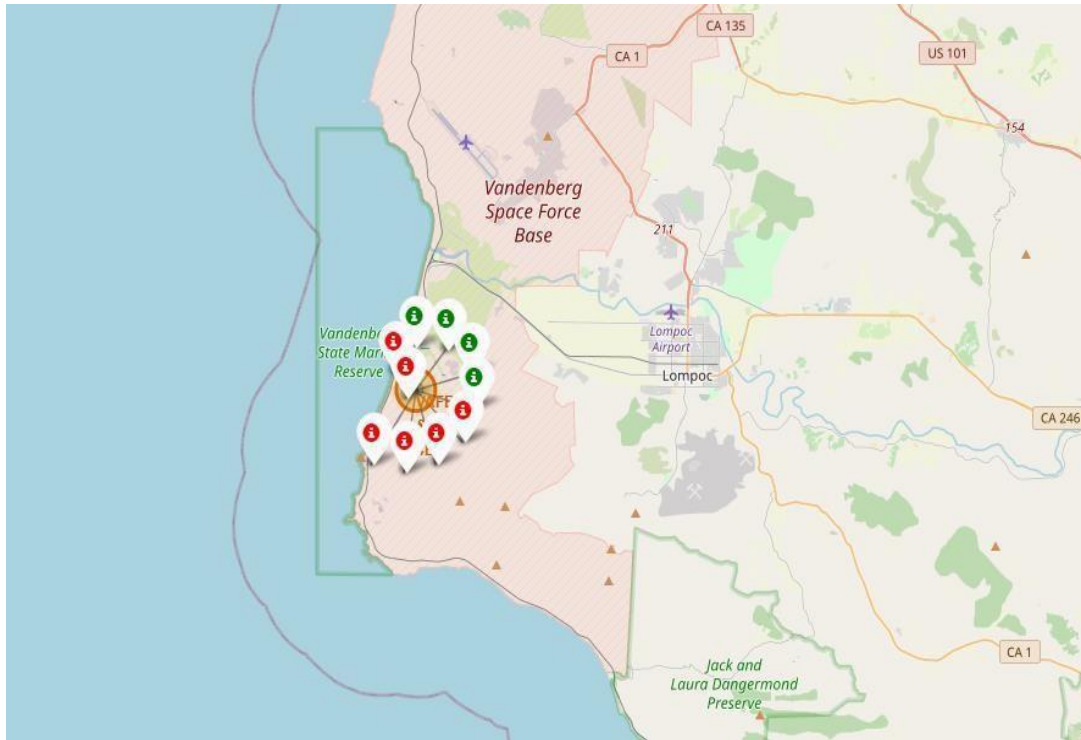
Build an interactive map with Folium

- 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.

Launch Site Locations



Color coded launch markers



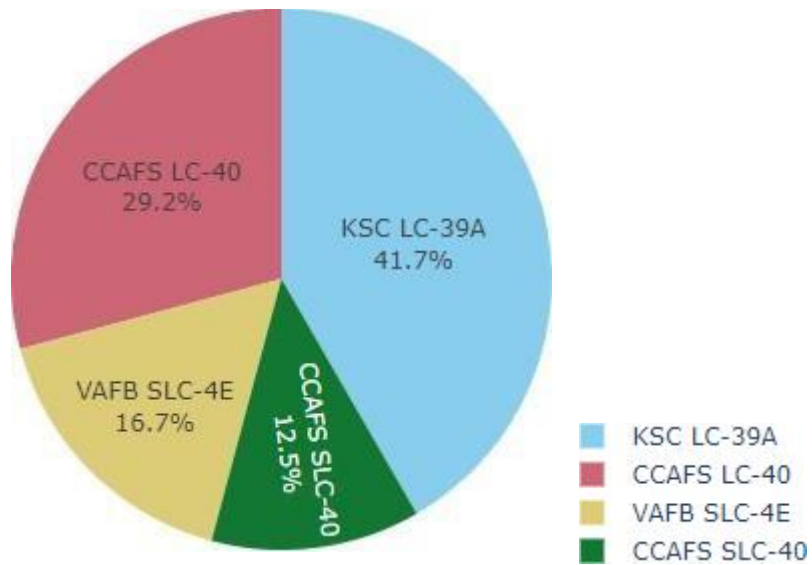
Using KSC LC-39A as an example, launch sites are very close to railways for large part and supply transportation. Launch sites are close to highways for human and supply transport. Launch sites are also close to coasts and relatively far from cities so that launch failures can land in the sea to avoid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



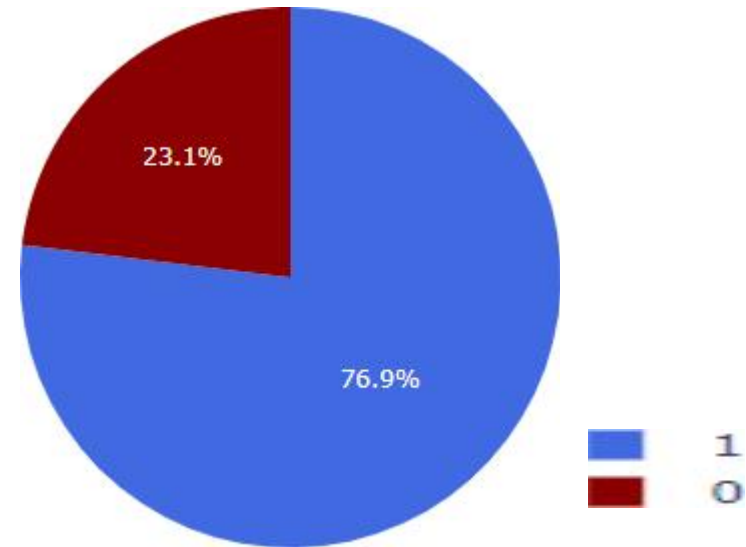
Build A Dashboard with Plotly Dash

- 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.

Successful launches across launch sites



Highest Success rate launch site

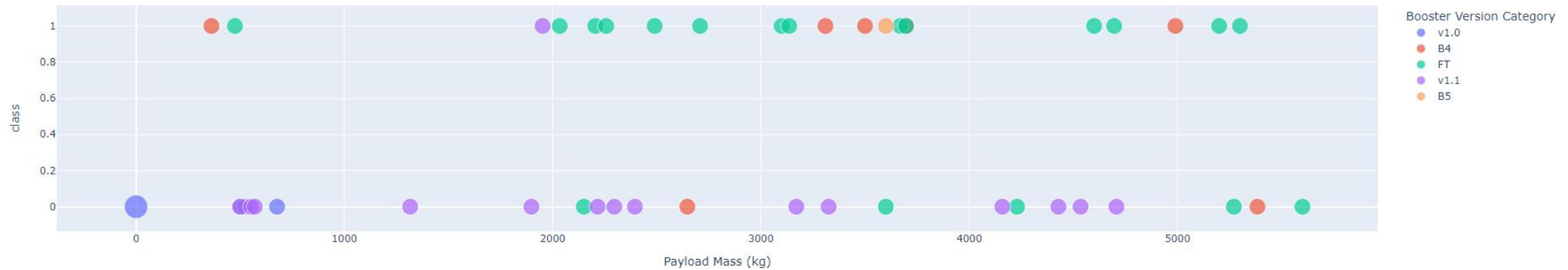


Payload Mass Vs Success Vs Booster version Category

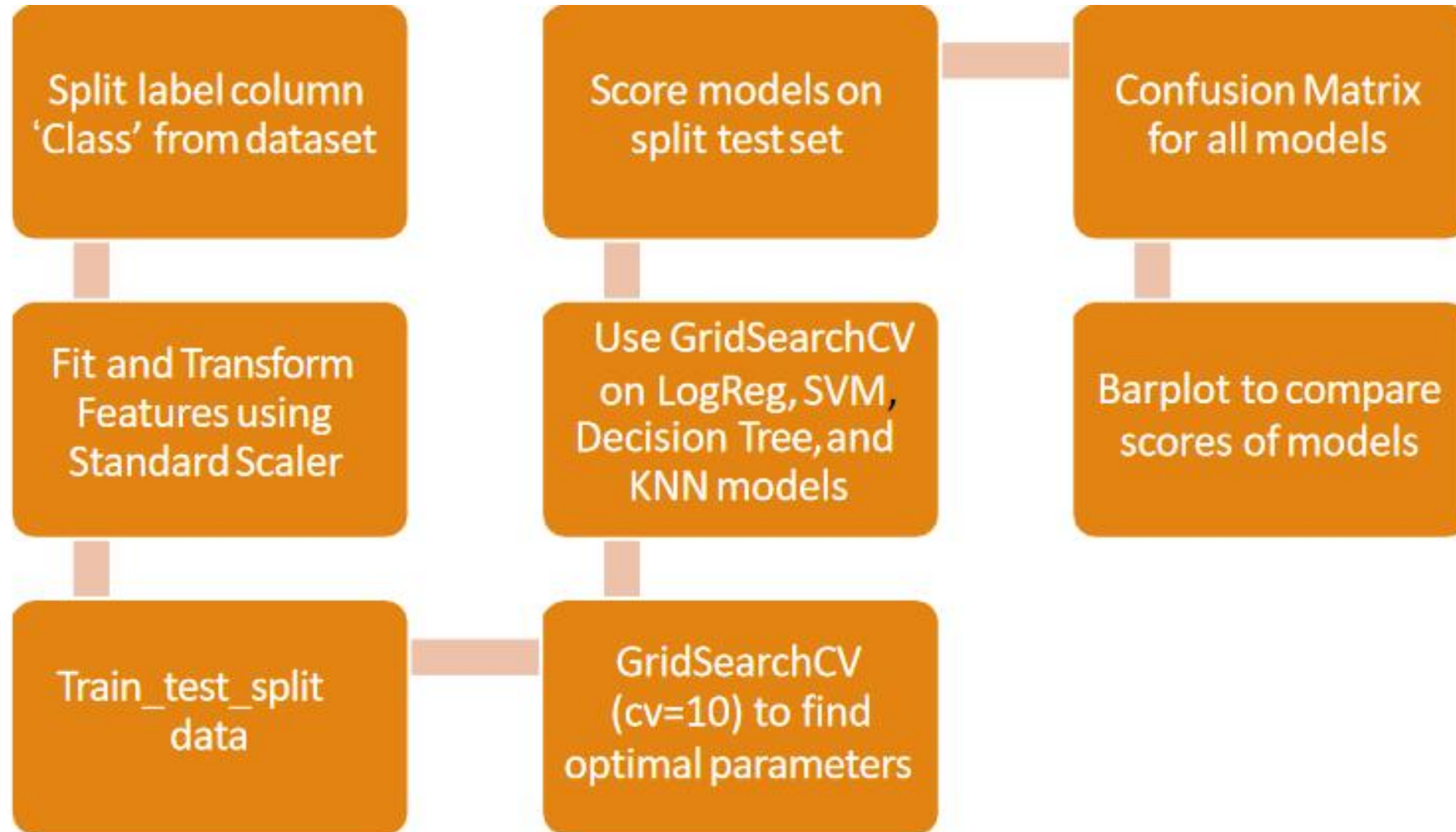
Payload range (Kg):



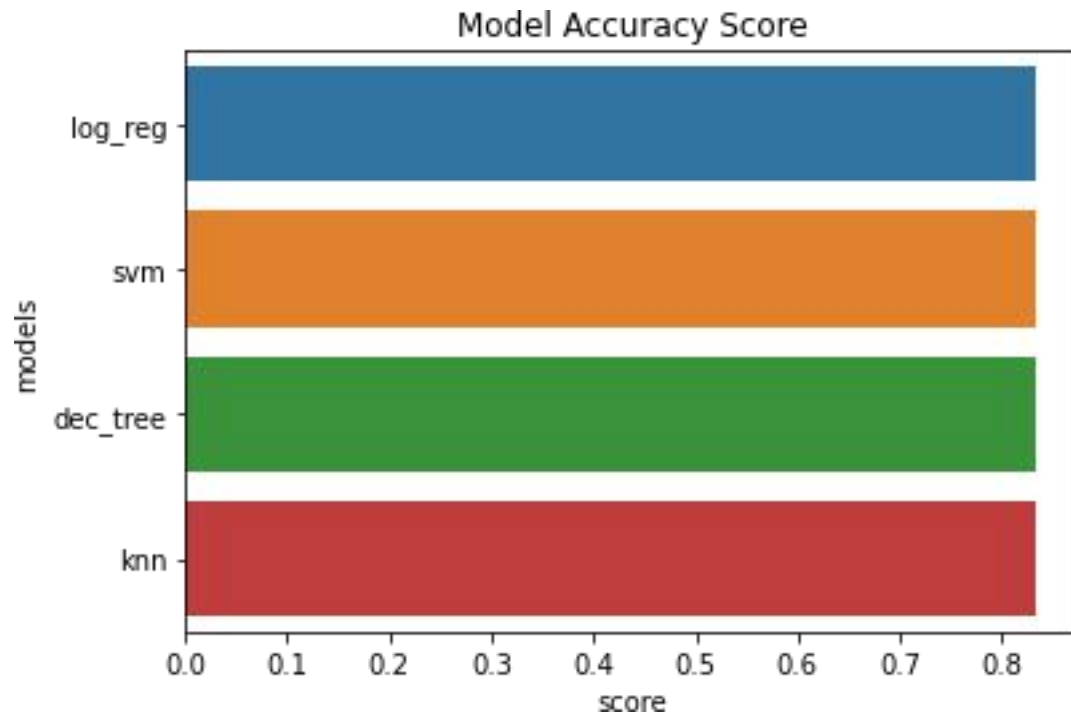
Payload Mass vs. Success vs. Booster Version Category



Predictive Analysis (Classification)

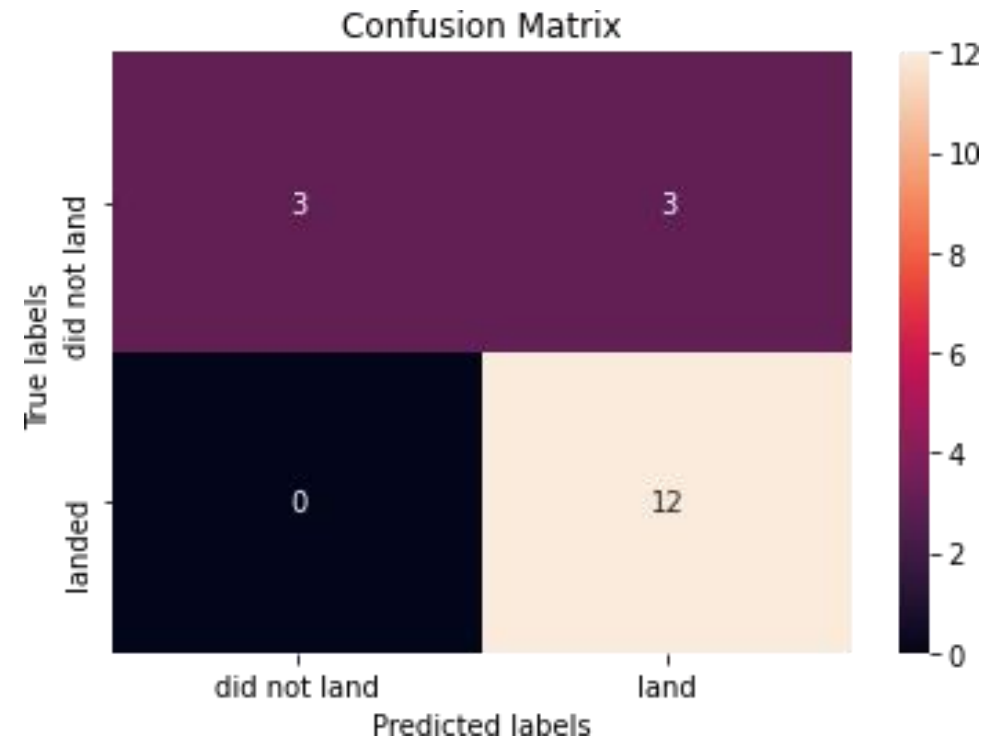


Classification Accuracy



- 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



- 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

We built a tool to help Space Y compete with SpaceX by predicting successful Stage 1 landings.

- This model analyzes data from a public SpaceX API and Wikipedia to predict if the first stage of a launch will land successfully.
- A successful landing saves Space Y an estimated \$100 million.
- The model uses data stored in a DB2 SQL database and includes a visualization dashboard for easy access.
- With an accuracy of around 83%, the model can help Space Y decide whether to proceed with a launch based on the predicted landing outcome.
- While the model is promising, collecting more data can further improve its accuracy and potentially identify the best machine learning approach.