

# **Data ScienceCapstone Project**

---

Saranya Suryadevara

3/14/2024

# Contents

---

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



## Executive Summary

Information gathered from the SpaceX Wikipedia page and public API. 'Class' is a newly created labels column that categorizes successful landings. Used dashboards, folium maps, SQL, visualization, and data exploration. Gathered pertinent columns to serve as the features. Used a single hot encoding to convert all categorical variables to binary. GridSearchCV was utilized to determine the optimal parameters for machine learning models using standardized data. Visualize each model's accuracy score.

## Introduction

---

- 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

---

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

A combination of web scraping data from a table in Space X's Wikipedia entry and API queries from the public API were used in the data collection procedure. The data collection flowchart from API will be displayed on the following slide, and the data collection flowchart from webscraping will be displayed on the one after that. FlightNumber, Date, BoosterVersion, PayloadMass, Orbit, LaunchSite, Result, Flights, GridFins, Reused, Legs, LandingPad, Block, ReusedCount, Serial, Longitude, Latitude are the Space X API Data Columns. Wikipedia Flight No., Launch location, Payload, PayloadMass, Orbit, Customer, Launch result, Version Booster, Booster landing, Date, and Time are the columns of webscrape data.

## 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

## 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



## EDA with SQL

---

- Loaded data set into IBM DB2Database.
- 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

## 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.

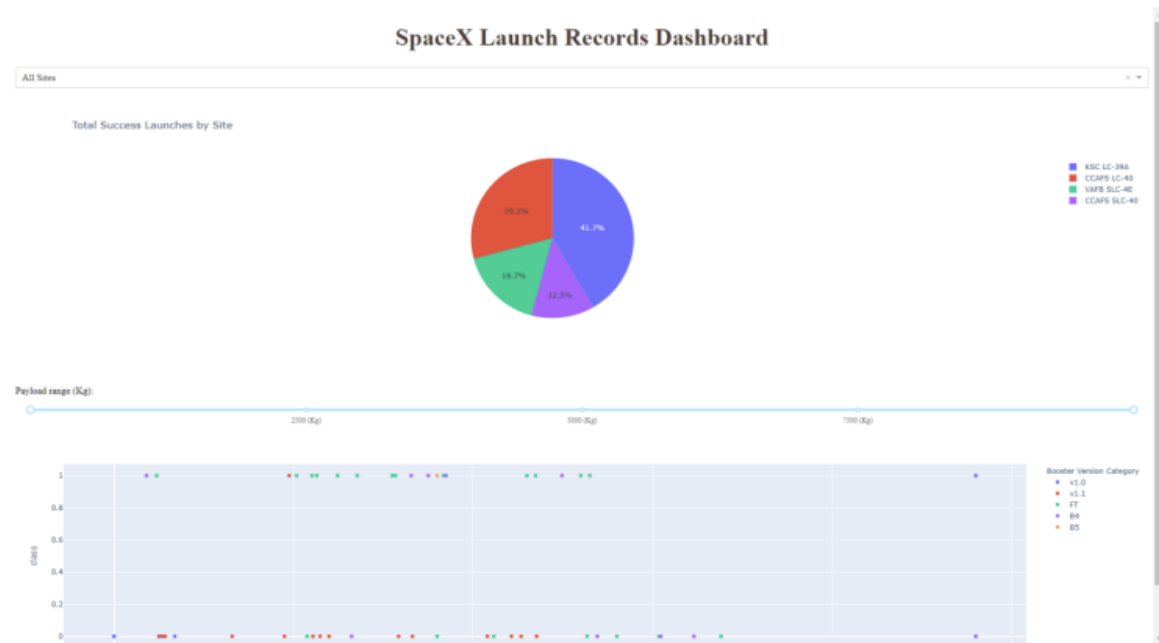
## Build a Dashboard with PlotlyDash

---

- 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.

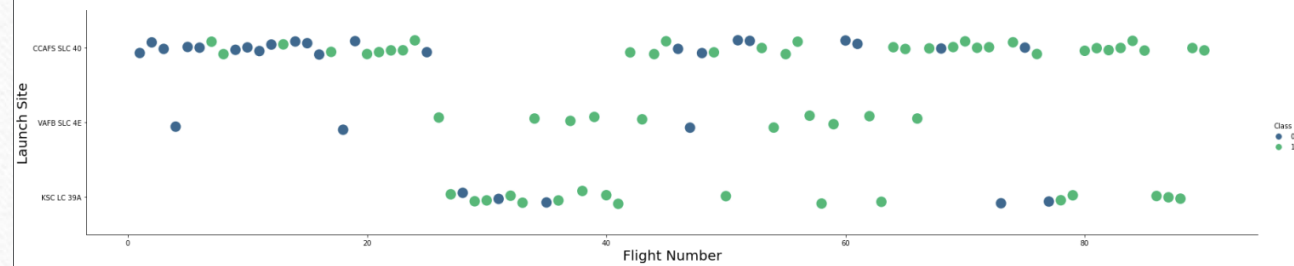


# Results



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.

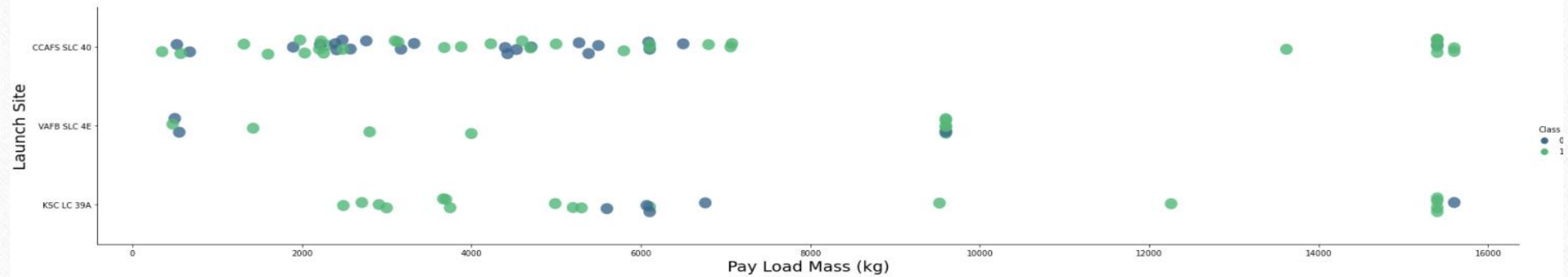
## Flight Number vs. Launch Site



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

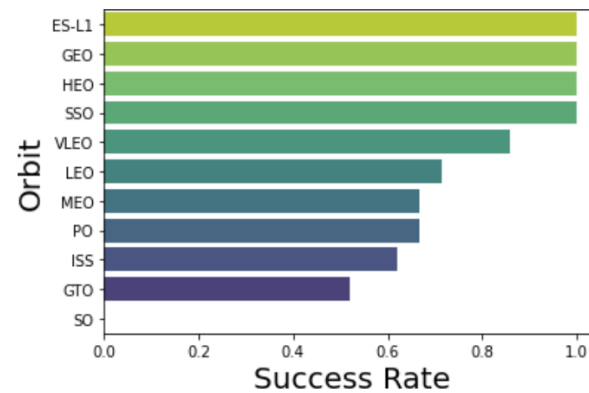


Green indicates successful launch; Purple indicates unsuccessful launch.

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



## Successrate vs. Orbittype



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)

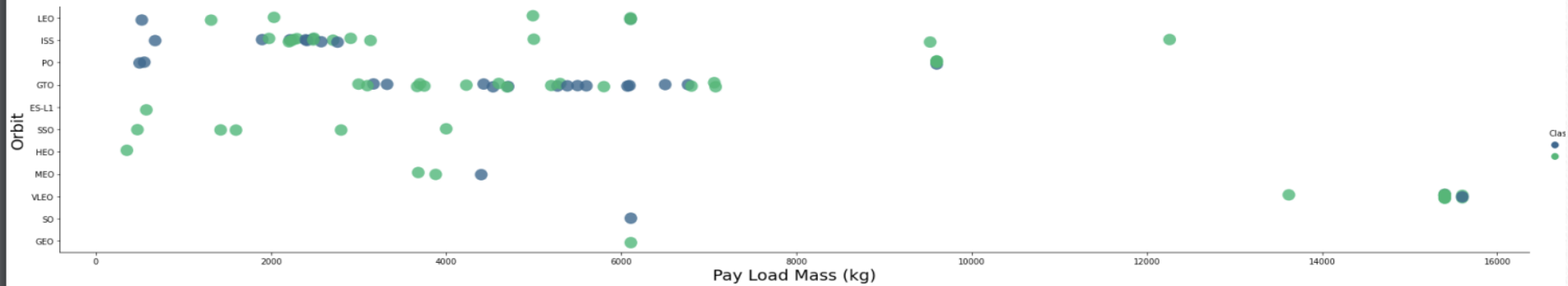
SSO (5) has 100% success rate

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

## Payload vs. Orbit type



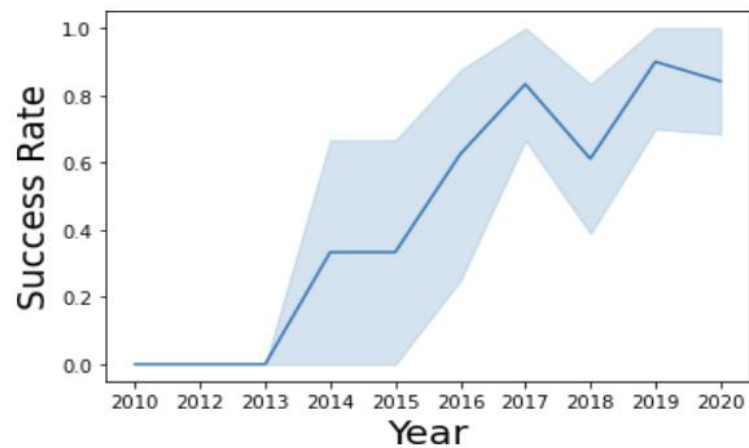
Green indicates successful launch; Purple indicates unsuccessful launch.

## 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

```
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

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

## Average Payload Mass by F9v1.1

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

```
* ibm_db_sa://ftb12020:***@0c77d6f2-5da9-48a9-81f8-80
Done.
```

avg_payload_mass_kg
---------------------

2928
------

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

## Average Payload Mass by F9v1.1

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

```
* ibm_db_sa://ftb12020:***@0c77d6f2-5da9-48a9-81f8-80
Done.
```

avg_payload_mass_kg
---------------------

2928
------

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 Pad Landing Date

```
%%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

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

---

```
%%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.bs2io90108kqb1od81cg.database
Done.
```

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

```
%%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

This 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.



## 2015 Failed Drone Ship Landing Records

```
%%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.bs21o90108kqb1od81cg.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

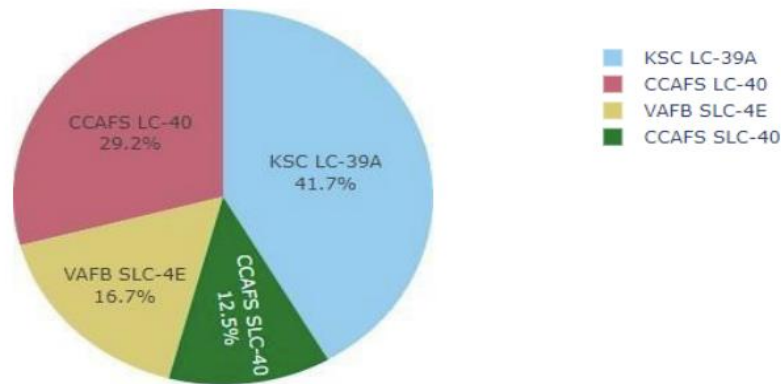
This query returns the Month, Landing Outcome, Booster Version, Payload Mass (kg), and Launch site of 2015 launches where stage 1 failed to land on a drone ship.

There were two such occurrences.

# Build a Dashboard with Plotly Dash

---

## Successful Launches Across Launch Sites

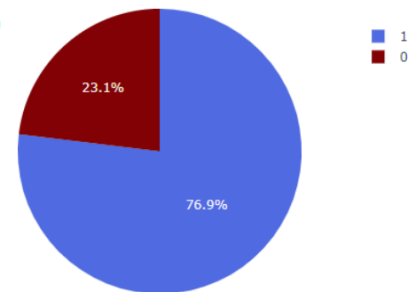


This is the distribution of successful landings across all launch sites. CCAFS LC-40 is the old name of CCAFS SLC-40 so CCAFS and KSC have the same amount of successful landings, but a majority of the successful landings were performed before the name change. VAFB has the smallest share of successful 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)



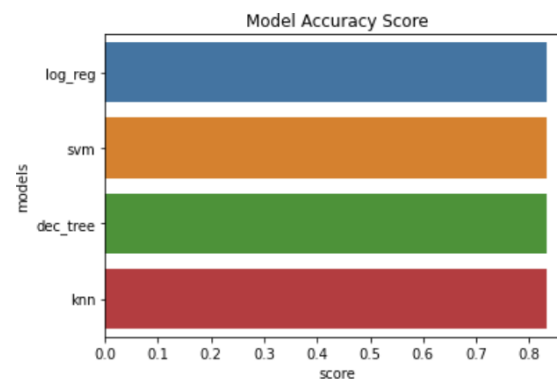
KSC LC-39A has the highest success rate with 10 successful landings and 3 failed landings.

## 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.

## Classification Accuracy



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

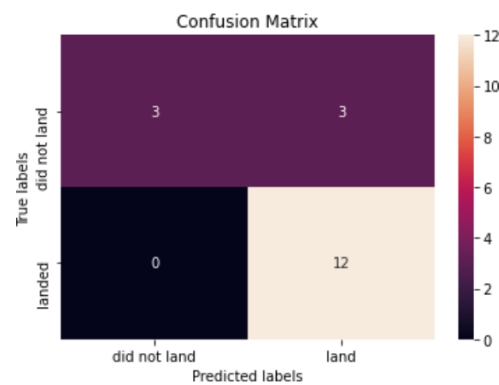
It should be noted that test size is small at only sample size of 18.

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.

## CONCLUSION

---

- 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