



# Winning Space Race with Data Science

Steven  
19/04/2024



# Outline

---

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

# Executive Summary

---

## **The following methodologies were used to analyze data:**

- Data Collection using web scraping and SpaceX API;
- Exploratory Data Analysis (EDA), including data wrangling, data visualization and interactive visual analytics;
- Machine Learning Prediction.

## **Summary of all results:**

- It was possible to collect valuable data from public sources;
- EDA allowed to identify which features are the best to predict success of launchings;
- Machine Learning Prediction showed the best model to predict which characteristics are
- important to drive this opportunity by the best way, using all collected data.

# Introduction

---

## Background

SpaceX, a leader in the space industry, strives to make space travel affordable for everyone. Its accomplishments include sending spacecraft to the international space station, launching a satellite constellation that provides internet access and sending manned missions to space. SpaceX can do this because the rocket launches are relatively inexpensive (\$62 million per launch) due to its novel reuse of the first stage of its Falcon 9 rocket. Other providers, which are not able to reuse the first stage, cost upwards of \$165 million each. By determining if the first stage will land, we can determine the price of the launch. To do this, we can use public data and machine learning models to predict whether SpaceX -or a competing company -can reuse the first stage.

## Explore

- How payload mass, launch site, number of flights, and orbits affect first-stage landing success
- Rate of successful landings over time
- Best predictive model for successful landing (binary classification)

Section 1

# Methodology

# Methodology

---

## Executive Summary

- Data collection methodology:
  - Data from Space X was obtained from 2 sources:
    - Space X API(<https://api.spacexdata.com/v4/rockets/>)
    - WebScraping  
([https://en.wikipedia.org/wiki/List\\_of\\_Falcon/ 9/ and Falcon Heavy launches](https://en.wikipedia.org/wiki/List_of_Falcon/9_and_Falcon_Heavy_launches))
- Perform data wrangling
  - Collected data was enriched by creating a landing outcome label based on outcome data after summarizing and analyzing features
- Perform exploratory data analysis (EDA) using visualization and SQL

# Methodology

## Executive Summary

- Perform interactive visual analytics using Folium and Plotly Dash
- Perform predictive analysis using classification models
  - Data that was collected until this step were normalized, divided in training and test data sets and evaluated by four different classification models, being the accuracy of each model evaluated using different combinations of parameters.

# Data Collection

---

Data sets were collected from

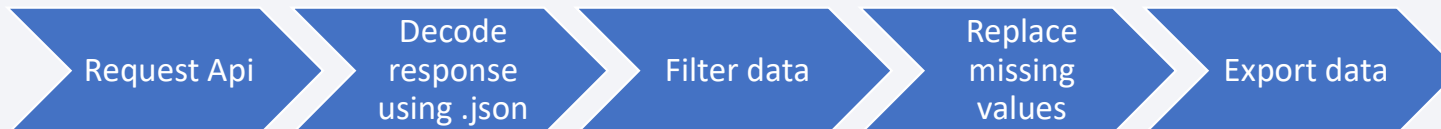
- Space X API (<https://api.spacexdata.com/v4/rockets/>)
- Wikipedia ([https://en.wikipedia.org/wiki/List\\_of\\_Falcon\\_9\\_and\\_Falcon\\_Heavy\\_launches](https://en.wikipedia.org/wiki/List_of_Falcon_9_and_Falcon_Heavy_launches)), using web scraping technics.



# Data Collection – SpaceX API

---

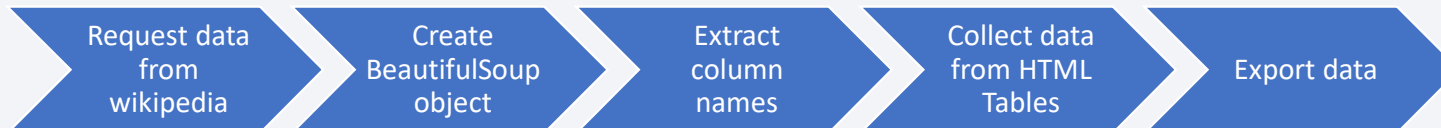
- SpaceX offers a public API from where data can be obtained and then used;
- This API was used according to the flowchart beside and then data is persisted.
- [https://github.com/StDierick/CourseraCapstone/blob/main/01\\_jupyter-labs-spacex-data-collection-api.ipynb](https://github.com/StDierick/CourseraCapstone/blob/main/01_jupyter-labs-spacex-data-collection-api.ipynb)



# Data Collection - Scraping

---

- Data from SpaceX launches can also be obtained from Wikipedia;
- Data are downloaded from Wikipedia according to the flowchart and then persisted.
- [https://github.com/StDierick/CourseraCapstone/blob/main/O2\\_jupyter-labs-webscraping.ipynb](https://github.com/StDierick/CourseraCapstone/blob/main/O2_jupyter-labs-webscraping.ipynb)



# Data Wrangling

---

- The .csv file from the first section contains the data that needed to be cleaned.
- The launch sites, orbit types and mission outcomes were cleaned up.
- The handful of mission outcome types were converted to a binary classification where 1 means that the Falcon 9 first stage landing was a success and 0 means that it was a failure.
- The new classification was added to the DataFrame for further analysis
- [https://github.com/StDierick/CourseraCapstone/blob/main/03\\_labs-jupyter-spacex-Data%20wrangling.ipynb](https://github.com/StDierick/CourseraCapstone/blob/main/03_labs-jupyter-spacex-Data%20wrangling.ipynb)

# EDA with Data Visualization

---

- **Charts**

- Flight Number vs. Payload
- Flight Number vs. Launch Site
- Payload Mass (kg) vs. Launch Site
- Payload Mass (kg) vs. Orbit type

- **Analysis**

- **View relationship** by using **scatter plots**. The variables could be useful for machine learning if a relationship exists
- **Show comparisons** among discrete categories with **bar charts**. Bar charts show the relationships among the categories and a measured value.

- [https://github.com/StDierick/CourseraCapstone/blob/main/05\\_jupyter-labs-eda-dataviz.ipynb](https://github.com/StDierick/CourseraCapstone/blob/main/05_jupyter-labs-eda-dataviz.ipynb)

# EDA with SQL

---

- The following SQL queries were performed:
  - Names of the unique launch sites in the space mission;
  - Top 5 launch sites whose name begin with the string 'CCA';
  - Total payload mass carried by boosters launched by NASA(CRS);
  - Average payload mass carried by booster version F9 v1.1;
  - Date when the first successful landing outcome in ground pad was achieved;
  - Names of the boosters which have success in drone ship and have payload mass between 4000 and 6000 kg;
  - Total number of successful and failure mission outcomes;
  - Names of the booster versions which have carried the maximum payload mass;
  - Failed landing outcomes in drone ship, their booster versions, and launch site names for in year 2015; and
  - Rank of the count of landing outcomes (such as Failure (drone ship) or Success (ground pad)) between the date 2010-06-04 and 2017-03-20.
- [https://github.com/StDierick/CourseraCapstone/blob/main/04\\_jupyter-labs-eda-sql-coursera\\_sqllite.ipynb](https://github.com/StDierick/CourseraCapstone/blob/main/04_jupyter-labs-eda-sql-coursera_sqllite.ipynb)

# Build an Interactive Map with Folium

---

- **Markers Indicating Launch Sites**

- Added **blue circle** at **NASA Johnson Space Center's coordinate** with a **popup label** showing its name using its latitude and longitude coordinates
- Added **red circles** at **all launch sites coordinates** with a **popup label** showing its name using its name using its latitude and longitude coordinates

- **Colored Markers of Launch Outcomes**

- Added **colored markers** of **successful (green)** and **unsuccessful (red) launches** at each launch site to show which launch sites have high success rates

- **Distances Between a Launch Site to Proximities**

- Added **colored lines** to **show distance between** launch site **CCAFS SLC-40 and** its proximity to the **nearest coastline, railway, highway, and city**
- [https://github.com/StDierick/CourseraCapstone/blob/main/6\\_1\\_lab\\_jupyter\\_launch\\_site\\_location-Folium%20lab.ipynb](https://github.com/StDierick/CourseraCapstone/blob/main/6_1_lab_jupyter_launch_site_location-Folium%20lab.ipynb)

# Build a Dashboard with Plotly Dash

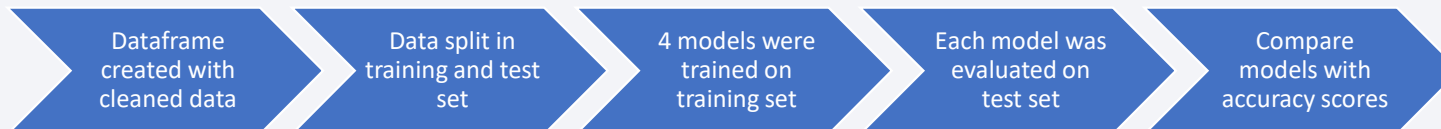
---

- **Dropdown List with Launch Sites**
  - Allow user to select all launch sites or a certain launch site
- **Slider of Payload Mass Range**
  - Allow user to select payload mass range
- **Pie Chart Showing Successful Launches**
  - Allow user to see successful and unsuccessful launches as a percent of the total
- **Scatter Chart Showing Payload Mass vs. Success Rate by Booster Version**
  - Allow user to see the correlation between Payload and Launch Success
- [https://github.com/StDierick/CourseraCapstone/blob/main/6\\_2\\_spaceX%20Dash%20App.ipynb](https://github.com/StDierick/CourseraCapstone/blob/main/6_2_spaceX%20Dash%20App.ipynb)

# Predictive Analysis (Classification)

---

- The dataset was split into training and testing sets.
- Logistic Regression, SVM (Support Vector Machine), Decision Tree, and KNN (k-Nearest Neighbors) machine learning models were trained on the training data set.
- Hyper-parameters were evaluated using GridSearchCV() and the best was selected using.
- Using the best hyper-parameters, each of the four models were scored on accuracy by using the testing data set.
- [https://github.com/StDierick/CourseraCapstone/blob/main/07\\_SpaceX\\_Machine\\_Learning\\_Prediction\\_Part\\_5.pyterlite.ipynb](https://github.com/StDierick/CourseraCapstone/blob/main/07_SpaceX_Machine_Learning_Prediction_Part_5.pyterlite.ipynb)





# Results

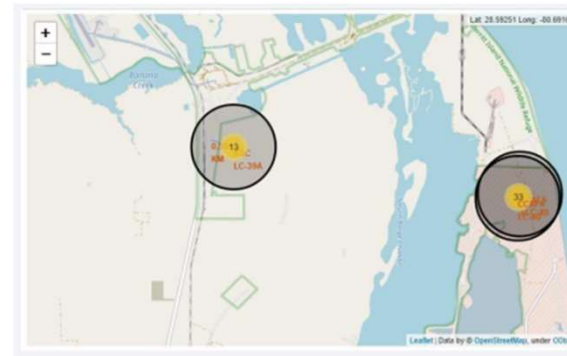
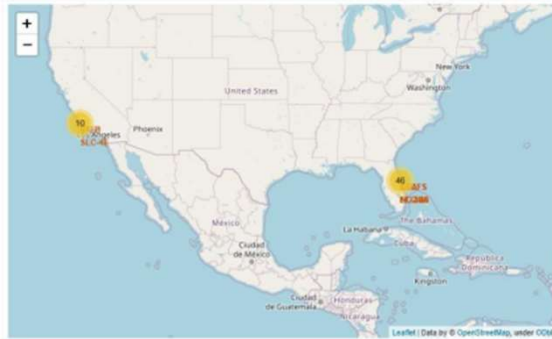
---

Exploratory data analysis results:

- Space X uses 4 different launch sites;
- The first launches were done to Space X itself and NASA;
- The average payload of F9 v1.1 booster is 2,928 kg;
- The first successful landing outcome happened in 2015 five years after the first launch;
- Many Falcon 9 booster versions were successful at landing in drone ships having payload above the average;
- Almost 100% of mission outcomes were successful;
- Two booster versions failed at landing in drone ships in 2015: F9 v1.1 B1012 and F9 v1.1 B1015;
- The number of landing outcomes became as better as years passed.

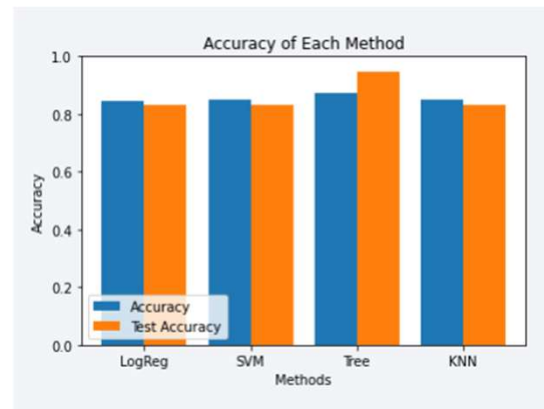
# Results

- Using interactive analytics was possible to identify that launch sites use to be in safety places, near sea, for example and have a good logistic infrastructure around.
- Most launches happens at east cost launch sites.



# Results

Predictive Analysis showed that Decision Tree Classifier is the best model to predict successful landings, having accuracy over 87% and accuracy for test data over 94%.





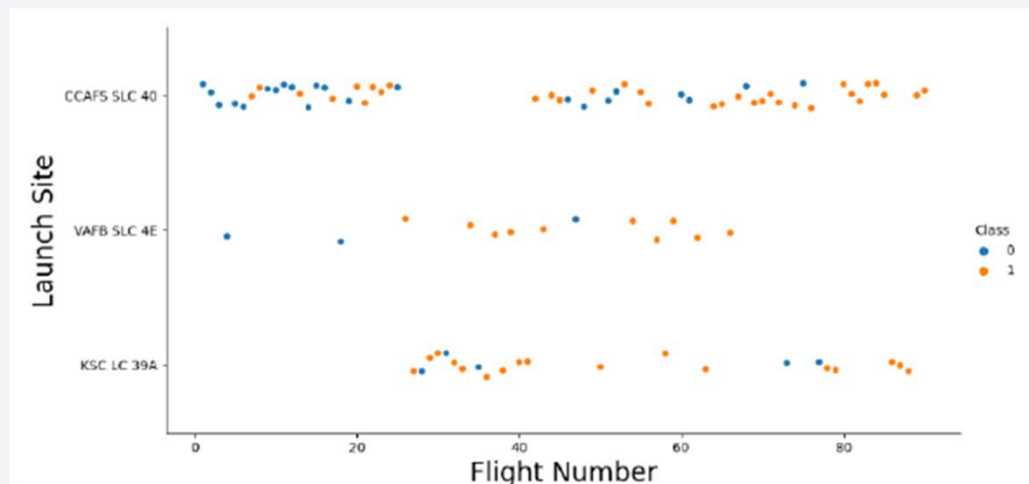
Section 2

# Insights drawn from EDA

# Flight Number vs. Launch Site

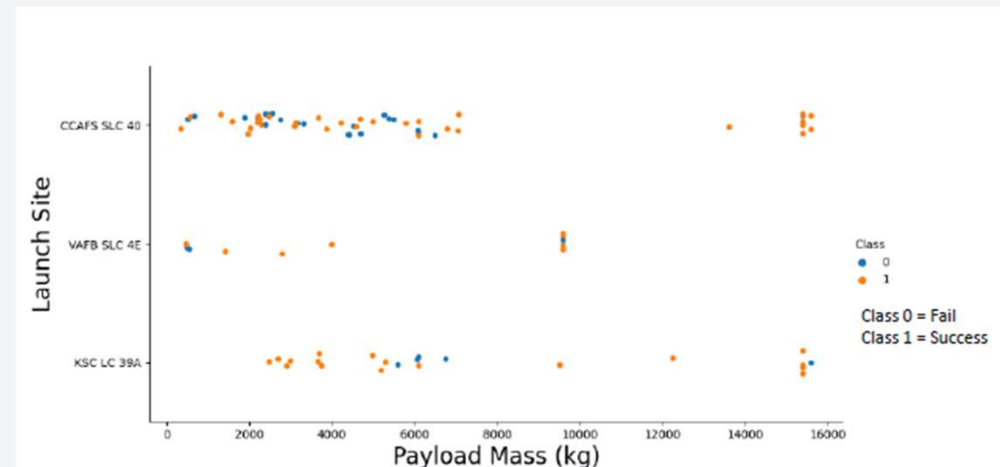
---

- According to the plot above, it's possible to verify that the best launch site nowadays is CCAF5 SLC 40, where most of recent launches were successful;
- In second place VAFB SLC 4E and third place KSC LC 39A;
- It's also possible to see that the general success rate improved over time.



# Payload vs. Launch Site

- Typically, the **higher** the **payload mass** (kg), the **higher** the **success rate**
- Most launches with a payload greater than 7,000 kg were successful
- KSC LC 39A has a 100% success rate for launches less than 5,500 kg
- VAFB SKC 4E has not launched anything greater than ~10,000 kg

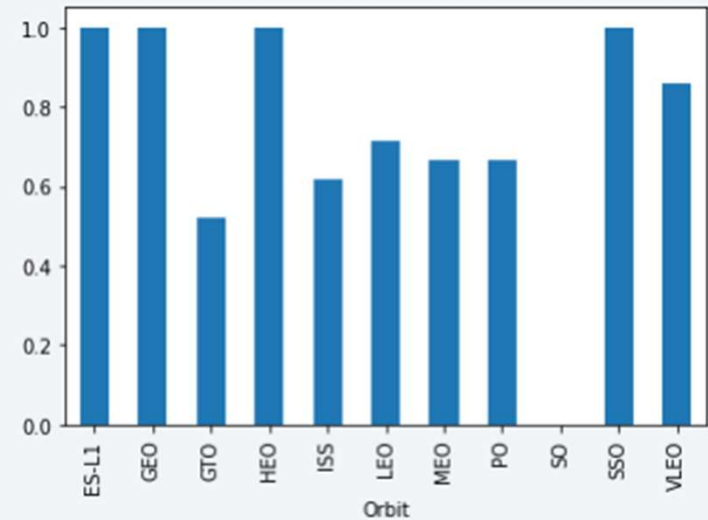




# Success Rate vs. Orbit Type

---

- The biggest success rates happens to orbits:
  - ES-L1;
  - GEO;
  - HEO;
  - SSO.
- Followed by:
  - VLEO (above 80%);
  - LFO (above 70%)



# Flight Number vs. Orbit Type

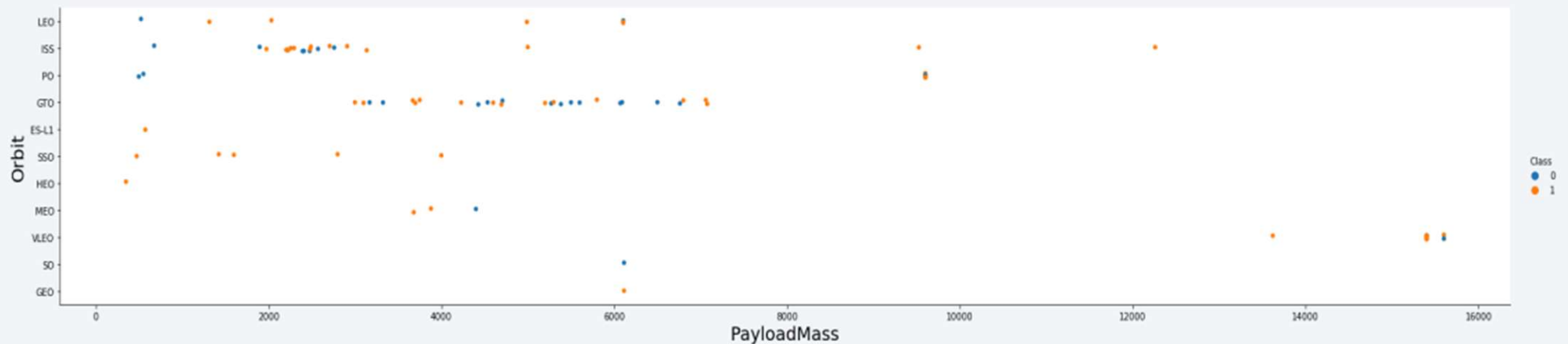
- Apparently, success rate improved over time to all orbits;
- VLEO orbit seems a new business opportunity, due to recent increase of its frequency





# Payload vs. Orbit Type

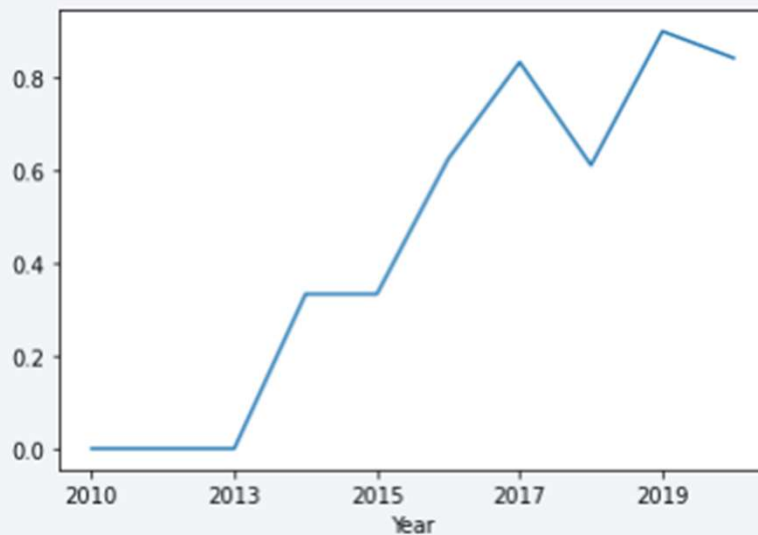
- Apparently, there is no relation between payload and success rate to orbit GTO;
- ISS orbit has the widest range of payload and a good rate of success;
- There are few launches to the orbits SO and GEO.



# Launch Success Yearly Trend

---

- Success rate started increasing in 2013 and kept until 2020;
- It seems that the first three years were a period of adjusts and improvement of technology.



# All Launch Site Names

---

- **Launch Site Names**

- CCAFS LC-40
- CCAFS SLC-40
- KSC LC-39A
- VAFB SLC-4<sup>E</sup>

- They are obtained by selecting unique occurrences of “launch\_site” values from the dataset.

# Launch Site Names Begin with 'CCA'

- Here we can see five samples of Cape Canaveral launches.

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

```
* ibm_db_sa://yyy33800:***@1bbf73c5-d84a-4bb0-85b9-ab1a4348f4a4.c3n41cmd0nqnk39u98g.databases.appdomain.cloud:32286/BLUDB
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	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

# Total Payload Mass

---

- **45,596 kg** (total) carried by boosters launched by NASA (CRS)

```
%sql SELECT SUM(PAYLOAD_MASS_KG_) \
      FROM SPACEXTBL \
      WHERE CUSTOMER = 'NASA (CRS)';

* ibm_db_sa://yyy33800:***@1bbf73c5-d84a-4l
sqlite:///my_data1.db
Done.

   1
---
45596
```

# Average Payload Mass by F9 v1.1

---

- **2,928 kg** (average) carried by booster version F9 v1.1

```
%sql SELECT AVG(PAYLOAD_MASS_KG_) \
FROM SPACEXTBL \
WHERE BOOSTER_VERSION = 'F9_v1.1';

* ibm_db_sa://yyy33800:***@1bbf73c5-d84a-4
sqlite:///my_data1.db
Done.

  1
---
2928
```

# First Successful Ground Landing Date

---

- First successful landing outcome on ground pad:

```
%sql SELECT MIN(DATE) \
FROM SPACEXTBL \
WHERE LANDING_OUTCOME = 'Success.(ground pad)'\
* ibm_db_sa://yyy33800:***@1bbf73c5-d84a-4bb0-85b0-...
sqlite:///my_data1.db
Done.

1
2015-12-22
```

- By filtering data by successful landing outcome on ground pad and getting the minimum value for date it's possible to identify the first occurrence, that happened on 12/22/2015.

## Successful Drone Ship Landing with Payload between 4000 and 6000

---

- Booster mass greater than 4,000 but less than 6,000
- JSCAT-14, JSCAT-16, SES-10, SES-11 / EchoStar 105

```
%sql SELECT PAYLOAD \
FROM SPACEXTBL \
WHERE LANDING_OUTCOME = 'Success (drone ship)' \
AND PAYLOAD_MASS_KG BETWEEN 4000 AND 6000;

* ibm_db_sa://yyy33800:***@1bbf73c5-d84a-4bb0-85b9-
sqlite:///my_data1.db
Done.
```

payload
JCSAT-14
JCSAT-16
SES-10
SES-11 / EchoStar 105



# Total Number of Successful and Failure Mission Outcomes

---

- 1 Failure in Flight
- 99 Success
- 1 Success (payload status unclear)

```
%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

---

- These are the boosters which have carried the maximum payload mass registered in the dataset.

```
%sql SELECT BOOSTER_VERSION \
FROM SPACEXTBL \
WHERE PAYLOAD_MASS_KG = (SELECT MAX(PAYLOAD_MASS_KG) FROM SPACEXTBL);

* 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

---

- Failed landing outcomes in drone ship, their booster versions, and launch site names for in year 2015

```
%sql SELECT substr(Date,4,2) as month, DATE,BOOSTER_VERSION, LAUNCH_SITE, [Landing _Outcome] \
FROM SPACEXTBL \
where [Landing _Outcome] = 'Failure (drone ship)' and substr(Date,7,4)='2015';
```

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

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

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

---

- Count of landing outcomes between 2010-06-04 and 2017-03-20 in descending order

```
%sql SELECT [Landing_Outcome], count(*) as count_outcomes \
FROM SPACEXTBL \
WHERE DATE between '04-06-2010' and '20-03-2017' group by [Landing_Outcome] order by count_outcomes DESC;
```

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

Done.

Landing_Outcome	count_outcomes
Success	20
No attempt	10
Success (drone ship)	8
Success (ground pad)	6
Failure (drone ship)	4
Failure	3
Controlled (ocean)	3
Failure (parachute)	2
No attempt	1

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 gradient on the left and a satellite photograph of Earth on the right. The Earth's surface is dark, with numerous bright yellow and orange lights representing city lights at night. The horizon of the Earth is visible, separating the dark surface from the deep blue of the sky.

Section 3

# Launch Sites Proximities Analysis

# All launch sites

- **Near Equator:** the closer the launch site to the equator, the **easier** it is **to launch** to equatorial orbit, and the more help you get from Earth's rotation for a prograde orbit. Rockets launched from sites near the equator get an **additional natural boost**-due to the rotational speed of earth -that **helps save the cost** of putting in extra fuel and boosters.



# Launch outcome by site

---

- Example of KSC LC-39A launch site launch outcomes

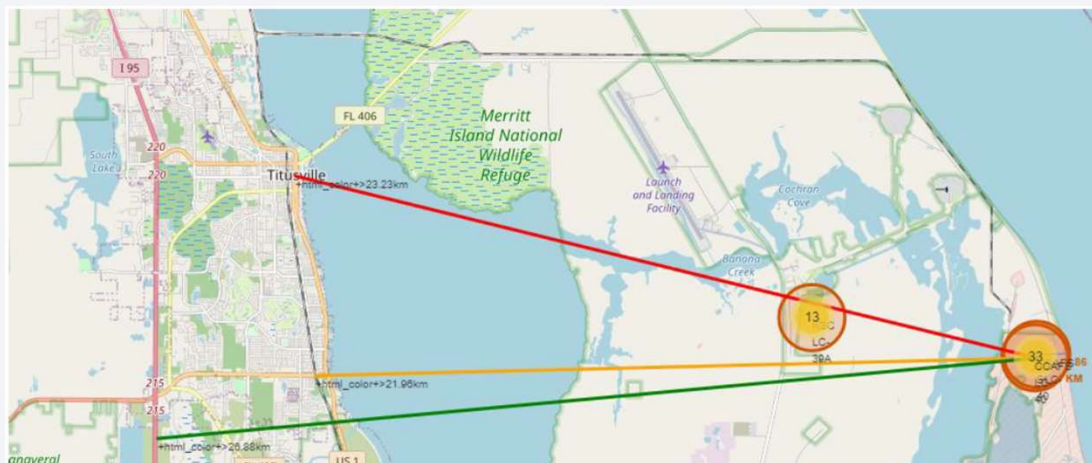


- Green markers indicate successful and red ones indicate failure

# Distance to proximities

## CCAFS SLC-40

- **0.86 km** from nearest coastline
- **21.96 km** from nearest railway
- **23.23 km** from nearest city
- **26.88 km** from nearest highway







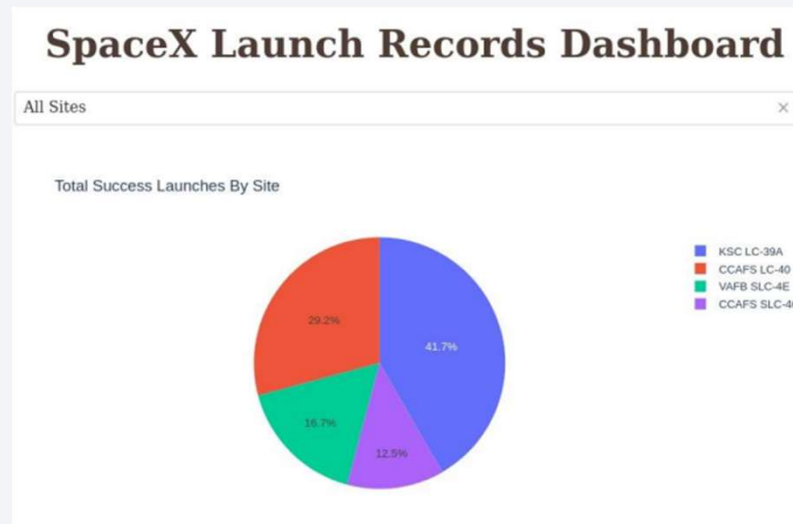
Section 4

# Build a Dashboard with Plotly Dash

# Successful launches by site

---

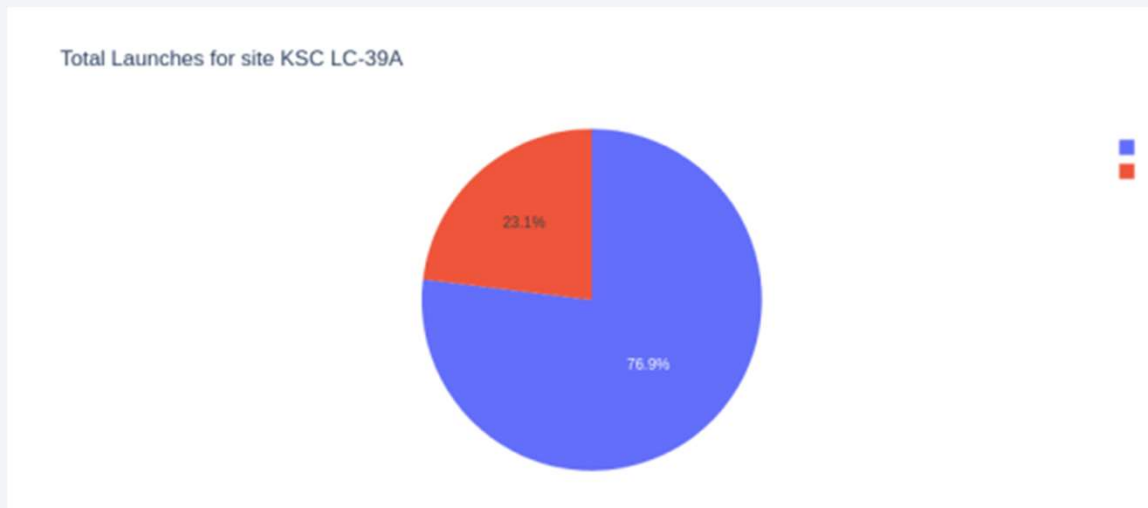
- The place from where launches are done seems to be a very important factor of success of missions.



# Launch Success Ratio for KSCLC-39A

---

- 76.9% of launches are successful in this site.



# Payload vs. Launch Outcome

- **Payloads between 2,000 kg and 5,000 kg have the highest success rate**





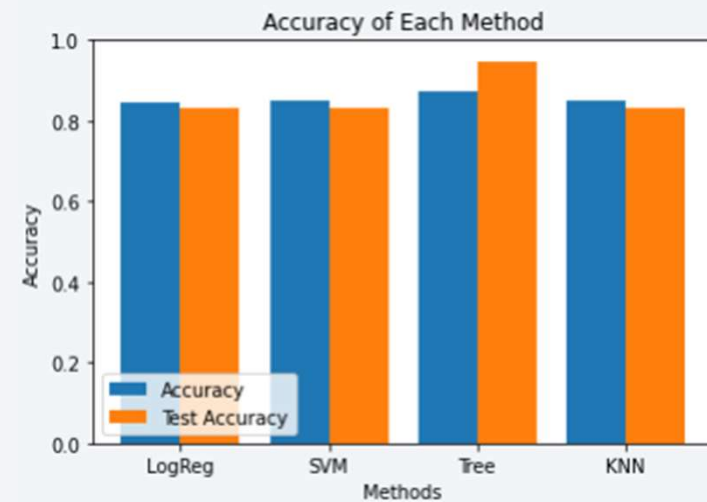
Section 5

# Predictive Analysis (Classification)

# Classification Accuracy

---

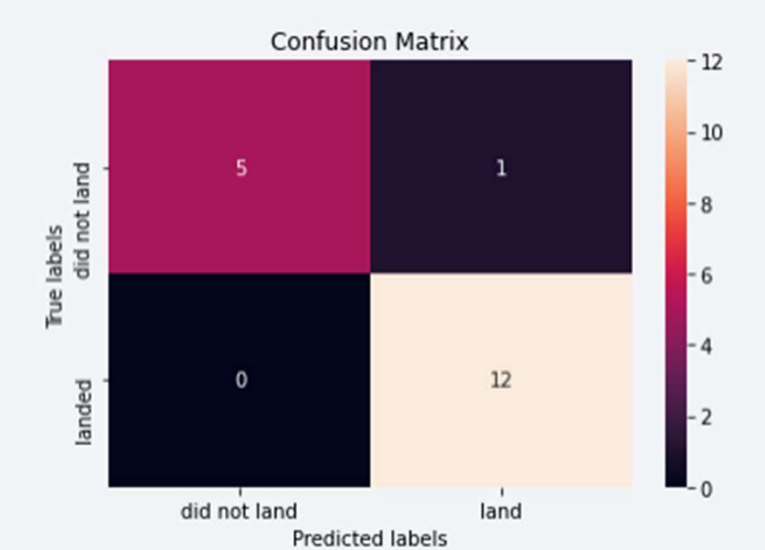
- Four classification models were tested, and their accuracies are plotted beside;
- The model with the highest classification accuracy is Decision Tree Classifier, which has accuracies over than 87%.



# Confusion Matrix

---

- Confusion matrix of Decision Tree Classifier proves its accuracy by showing the big numbers of true positive and true negative compared to the false ones.



# Conclusions

---

- Different data sources were analyzed, refining conclusions along the process;
- The best launch site is KSCLC-39A;
- Launches above 7,000kg are less risky;
- Although most of mission outcomes are successful, successful landing outcomes seem to improve over time, according the evolution of processes and rockets;
- Decision Tree Classifier can be used to predict successful landings and increase profits.



# Appendix

---

- All coding and datasets can be found on GitHub

<https://github.com/StDierick/CourseraCapstone/tree/main>

Thank you!

