# Outline

- Executive Summary

- Introduction

- Methodology

- Results

- Conclusion

- Appendix

# Executive Summary

- Summary of methodologies

  - Data Collection with Web Scraping and API intergration

  - Data Wrangling

  - Exploratory Analysis Using SQL and Data Visualization

  - Interactive Visual Analytics and Dashboard using Ploty and Dash

  - Predictive Analysis – SVM, KNN, Decision tree, Logistic Regression

- Summary of all results

  - Different launch sites have different success rate. CCAFS LC-40, has a success rate of 60 %, while KSC LC-39A and VAFB SLC 4E has a success rate of 77%.

  - there are no  rockets  launched for  heavypayload mass (10~k kg) in VAFB-SLC  launch site

  - All model using different methods presents similar accuracy. Succuss rate increased over time and remain steady in recent years

# Introduction

- Project background

  - Space Y (the new company) that would like to compete with SpaceX founded by Billionaire industrialist Allon Musk. Your job is to determine the price of each launch. You will do this by gathering information about Space X and creating dashboards for your team. You will also determine if SpaceX will reuse the first stage. Instead of using rocket science to determine if the first stage will land successfully, you will train a machine learning model and use public information to predict if SpaceX will reuse the first stage.

- Problems you want to find answers

  - What determines a successful landing?

  - What are the factors?

Section 1

# Methodology

# Methodology

Executive Summary

- Data collection methodology:

  - Request to the SpaceX API

- Perform data wrangling

  - Exploratory Data analysis applied to gain insight of data set

  - Training Labels assigned and categorical data encoded with one-hot

- 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

- Describe how data sets were collected.

    - Use HTTP REQUESTS to access SpaceX API and get a json file

    - Covert retrieved json file to dataframe use json_normalize().

    - Filter NaN and missing value.

    - Scrap additional information from Wikipedia about Falcon-9 historical launches. This was done with beautiful soup library.

    - Tables were parsed and converted into a readable pandas dataframe.

# Data Collection – SpaceX API

- Data download through SpaceX API and assigned to a dataframe for easy manipulation

- https://github.com/td121/Applied-Data-Science-Capstone/blob/main/jupyter-labs-spacex-data-collection-api.ipynb

# Data Collection - Scraping

- https://github.com/td121/Ap
  plied-Data-Science-
  Capstone/blob/main/jupyter-
  labs-webscraping.ipynb

# Data Wrangling

- Exploratory Data Analysis and Determine Training Labels
- We calculate and update the table:
  - 1. the number of launches on each site
  - 2. occurrence of each orbit
  - 3. number of mission outcome
  - 4. Add new landing outcome label

https://github.com/td121/Applied-Data-Science-Capstone/blob/main/labs-jupyter-spacex-Data%20wrangling.ipynb

# EDA with Data Visualization

- Fig.1 : the relationship between success rate of each orbit type

- Fig.2 : the relationship between Payload and Launch Site

- https://github.com/td121/Applied-Data-Science-Capstone/blob/main/jupyter-labs-eda-dataviz.ipynb

# EDA with SQL

- Load the SpaceX dataset into a PostgreSQL database.

- Apply EDA with SQL to get insight from the data. Following are considered:

    - Names of unique launch

    - Total payload mass carried by boosters launched by NASA (CRS)

    - Mean payload mass carried by booster version F9 v1.1

    - Total number of successful and failure mission outcomes

    - Total number of successful and failure mission outcomes

# Build an Interactive Map with Folium

- Mark all launch sites on a map

- Mark the success and failed launches for each site on the map

- Distance between a launch site to its proximities


- https://github.com/td121/Applied-Data-Science-Capstone/blob/main/lab_jupyter_launch_site_location.ipynb

# Build a Dashboard with Plotly Dash

- A dropdown list consist of a summarized item and each launch site

- A pie chart to show successful launches with respect to different sites

- A slider to adjust payload size

- Scatter plot show correlation between payload size and successful rate



https://github.com/td121/Applied-Data-Science-Capstone/blob/main/jupyter-plotly-dash.ipynb

14

# Predictive Analysis (Classification)

- Find the method performs best using test data

  - To Numpy array conversion of successful rate

  - Data standardization and split training data

  - Method performed

    - Logistic regression

    - Support vector machines (SVMs)

    - Decision tree

    - k-nearest neighbours algorithm (KNNs)

- https://github.com/td121/Applied-Data-Science-Capstone/blob/main/SpaceX_Machine%20Learning%20Prediction_Part_5.ipynb

# Results

- Exploratory data analysis results

- Interactive analytics demo in screenshots

- Predictive analysis results

Section 2

# Insights drawn from EDA

# Flight Number vs. Launch Site



WE can see more launches taken place in VAFB SLC 4E are successful. CCAFS SLC 40 presents a diverse range of success/fail rate.

# Payload vs. Launch Site



Pay Load Mass (kg) VS the launch site

- Payload Vs. Launch Site scatter point chart show that the VAFB-SLC launch site there are no rockets launched for heavy payload mass(greater than 10000).

# Success Rate vs. Orbit Type

# Flight Number vs. Orbit Type



It can be seen that 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



With heavy payloads the successful landing or positive landing rate are more for Polar,LEO and ISS.

# Launch Success Yearly Trend



Year VS Success rate

It is observe that the sucess rate since 2013 kept increasing till 2020

# All Launch Site Names

```python
# Select relevant sub-columns: `Launch Site`, `Lat(Latitude)`, `Long(Longitude)`, `class`
spacex_df = spacex_df[['Launch Site', 'Lat', 'Long', 'class']]
launch_sites_df = spacex_df.groupby(['Launch Site'], as_index=False).first()
launch_sites_df = launch_sites_df[['Launch Site', 'Lat', 'Long']]
launch_sites_df
```

|   | Launch Site | Lat | Long |
|---|---|---|---|
| 0 | CCAFS LC-40 | 28.562302 | -80.577356 |
| 1 | CCAFS SLC-40 | 28.563197 | -80.576820 |
| 2 | KSC LC-39A | 28.573255 | -80.646895 |
| 3 | VAFB SLC-4E | 34.632834 | -120.610745 |

# Launch Site Names Begin with 'CCA'



Display 5 records where launch sites begin with the string 'CCA'

```
task_2 = '''
        SELECT *
        FROM SpaceX
        WHERE LaunchSite LIKE 'CCA%'
        LIMIT 5
        '''
create_pandas_df(task_2, database=conn)
```

| | date | time | boosterversion | launchsite | payload | payloadmasskg | orbit | customer | missionoutcome | landingoutcome |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 2010-04-06 | 18:45:00 | F9 v1.0 B0003 | CCAFS LC-40 | Dragon Spacecraft Qualification Unit | 0 | LEO | SpaceX | Success | Failure (parachute) |
| 1 | 2010-08-12 | 15:43:00 | F9 v1.0 B0004 | CCAFS LC-40 | Dragon demo flight C1, two CubeSats, barrel of... | 0 | LEO (ISS) | NASA (COTS) NRO | Success | Failure (parachute) |
| 2 | 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 |
| 3 | 2012-08-10 | 00:35:00 | F9 v1.0 B0006 | CCAFS LC-40 | SpaceX CRS-1 | 500 | LEO (ISS) | NASA (CRS) | Success | No attempt |
| 4 | 2013-01-03 | 15:10:00 | F9 v1.0 B0007 | CCAFS LC-40 | SpaceX CRS-2 | 677 | LEO (ISS) | NASA (CRS) | Success | No attempt |

# Total Payload Mass

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

```
task_3 = '''
        SELECT SUM(PayloadMassKG) AS Total_PayloadMass
        FROM SpaceX
        WHERE Customer LIKE 'NASA (CRS)'
        '''
create_pandas_df(task_3, database=conn)
```

| | total_payloadmass |
|---|---|
| 0 | 45596 |

# Average Payload Mass by F9 v1.1

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

```
task_4 = '''
        SELECT AVG(PayloadMassKG) AS Avg_PayloadMass
        FROM SpaceX
        WHERE BoosterVersion = 'F9 v1.1'
        '''
create_pandas_df(task_4, database=conn)
```

|   | avg_payloadmass |
|---|-----------------|
| 0 | 2928.4          |

# First Successful Ground Landing Date

```
task_5 = '''
        SELECT MIN(Date) AS FirstSuccessfull_landing_date
        FROM SpaceX
        WHERE LandingOutcome LIKE 'Success (ground pad)'
        '''
create_pandas_df(task_5, database=conn)
```

| | firstsuccessfull_landing_date |
|---|---|
| 0 | 2015-12-22 |

# Successful Drone Ship Landing with Payload between 4000 and 6000

```
task_6 = '''
        SELECT BoosterVersion
        FROM SpaceX
        WHERE LandingOutcome = 'Success (drone ship)'
            AND PayloadMassKG > 4000
            AND PayloadMassKG < 6000
        '''

create_pandas_df(task_6, database=conn)
```

|   | boosterversion |
|---|---|
| 0 | F9 FT B1022 |
| 1 | F9 FT B1026 |
| 2 | F9 FT B1021.2 |
| 3 | F9 FT B1031.2 |

# Total Number of Successful and Failure Mission Outcomes

**List the total number of successful and failure mission outcomes**

```python
task_7a = '''
        SELECT COUNT(MissionOutcome) AS SuccessOutcome
        FROM SpaceX
        WHERE MissionOutcome LIKE 'Success%'
        '''

task_7b = '''
        SELECT COUNT(MissionOutcome) AS FailureOutcome
        FROM SpaceX
        WHERE MissionOutcome LIKE 'Failure%'
        '''
print('The total number of successful mission outcome is:')
display(create_pandas_df(task_7a, database=conn))
print()
print('The total number of failed mission outcome is:')
create_pandas_df(task_7b, database=conn)
```

The total number of successful mission outcome is:

| | successoutcome |
|---|---|
| 0 | 100 |

The total number of failed mission outcome is:

| | failureoutcome |
|---|---|
| 0 | 1 |

# Boosters Carried Maximum Payload

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

```
task_8 = '''
        SELECT BoosterVersion, PayloadMassKG
        FROM SpaceX
        WHERE PayloadMassKG = (
                              SELECT MAX(PayloadMassKG)
                              FROM SpaceX
                              )
        ORDER BY BoosterVersion
        '''
create_pandas_df(task_8, database=conn)
```

| | boosterversion | payloadmasskg |
|---|---|---|
| 0 | F9 B5 B1048.4 | 15600 |
| 1 | F9 B5 B1048.5 | 15600 |
| 2 | F9 B5 B1049.4 | 15600 |
| 3 | F9 B5 B1049.5 | 15600 |
| 4 | F9 B5 B1049.7 | 15600 |
| 5 | F9 B5 B1051.3 | 15600 |
| 6 | F9 B5 B1051.4 | 15600 |
| 7 | F9 B5 B1051.6 | 15600 |
| 8 | F9 B5 B1056.4 | 15600 |
| 9 | F9 B5 B1058.3 | 15600 |
| 10 | F9 B5 B1060.2 | 15600 |
| 11 | F9 B5 B1060.3 | 15600 |

# 2015 Launch Records

List the failed landing_outcomes in drone ship, their booster versions, and launch site names for in year 2015

```
task_9 = '''
        SELECT BoosterVersion, LaunchSite, LandingOutcome
        FROM SpaceX
        WHERE LandingOutcome LIKE 'Failure (drone ship)'
            AND Date BETWEEN '2015-01-01' AND '2015-12-31'
        '''
create_pandas_df(task_9, database=conn)
```

|   | boosterversion | launchsite | landingoutcome |
|---|----------------|------------|----------------|
| 0 | F9 v1.1 B1012 | CCAFS LC-40 | Failure (drone ship) |
| 1 | F9 v1.1 B1015 | CCAFS LC-40 | Failure (drone ship) |

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

Rank the count of landing outcomes (such as Failure (drone ship) or Success (ground pad))

```
task_10 = '''
        SELECT LandingOutcome, COUNT(LandingOutcome)
        FROM SpaceX
        WHERE DATE BETWEEN '2010-06-04' AND '2017-03-20'
        GROUP BY LandingOutcome
        ORDER BY COUNT(LandingOutcome) DESC
        '''
create_pandas_df(task_10, database=conn)
```

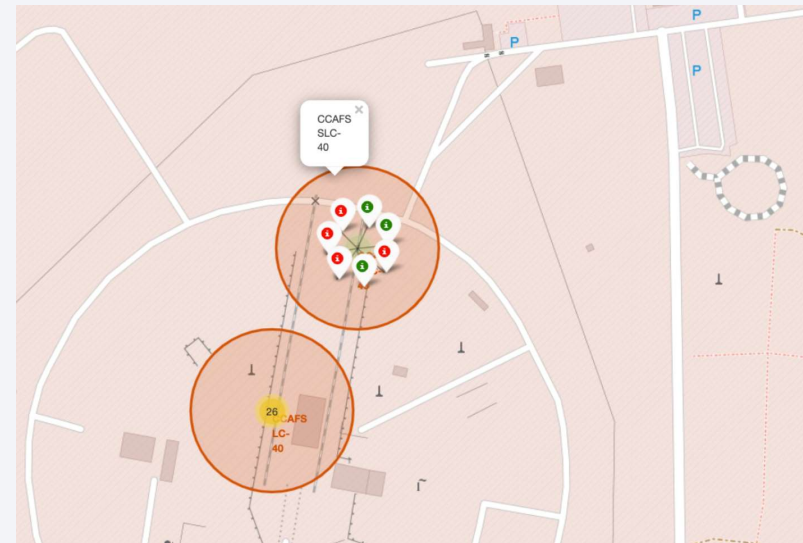|   | landingoutcome | count |
|---|----------------|-------|
| 0 | No attempt | 10 |
| 1 | Success (drone ship) | 6 |
| 2 | Failure (drone ship) | 5 |
| 3 | Success (ground pad) | 5 |
| 4 | Controlled (ocean) | 3 |
| 5 | Uncontrolled (ocean) | 2 |
| 6 | Precluded (drone ship) | 1 |
| 7 | Failure (parachute) | 1 |

Section 3

# Launch Sites
# Proximities Analysis
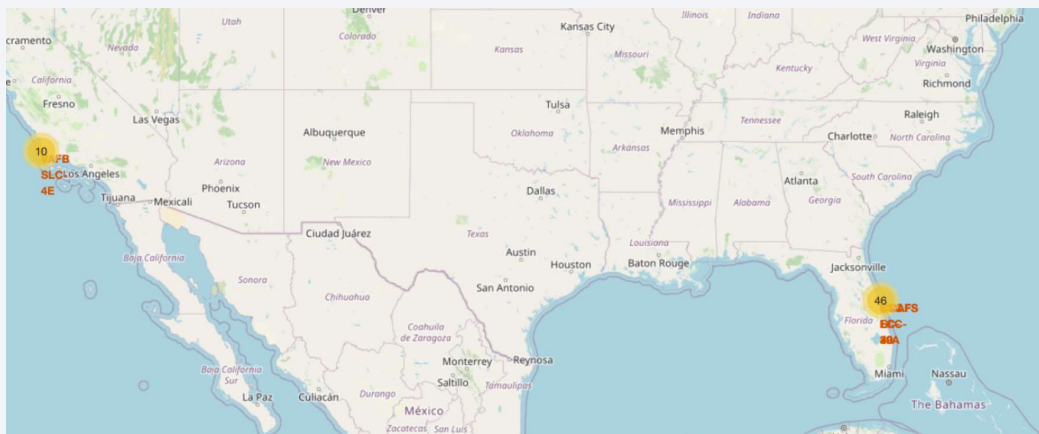
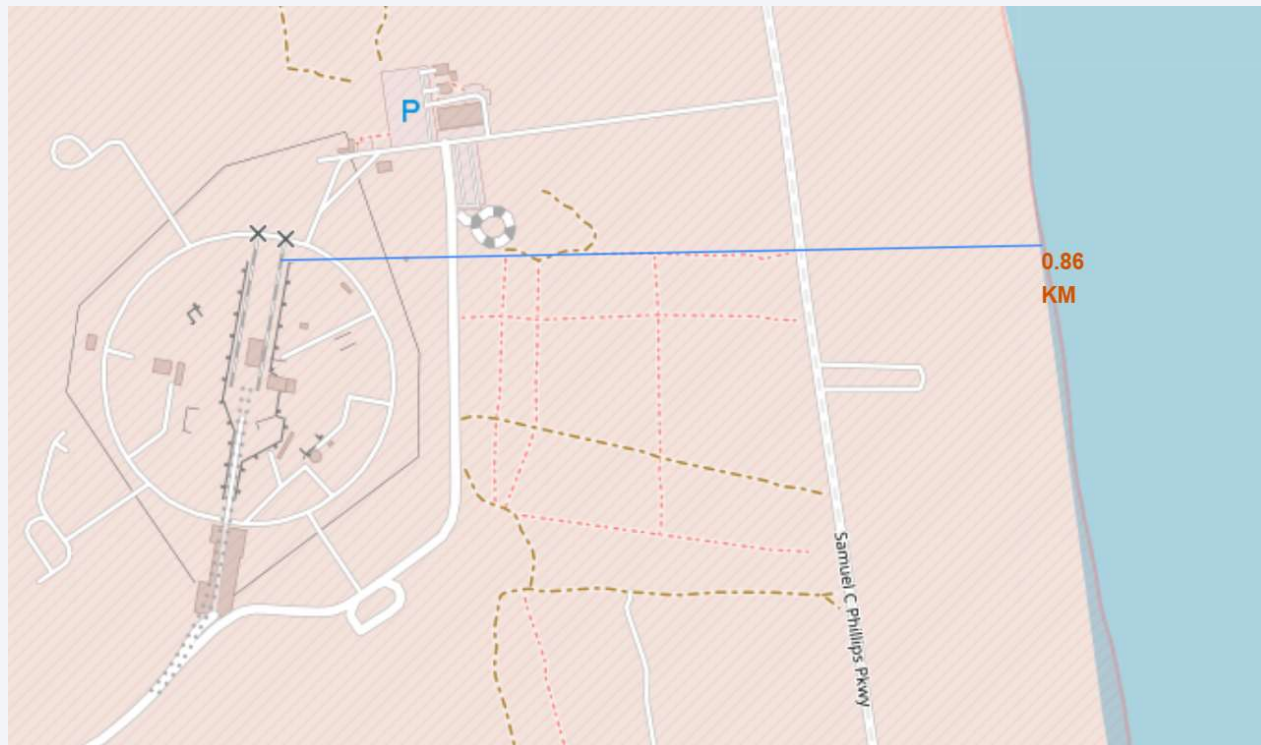# <Folium Map Screenshot 1>

- Launch site's locations on the map

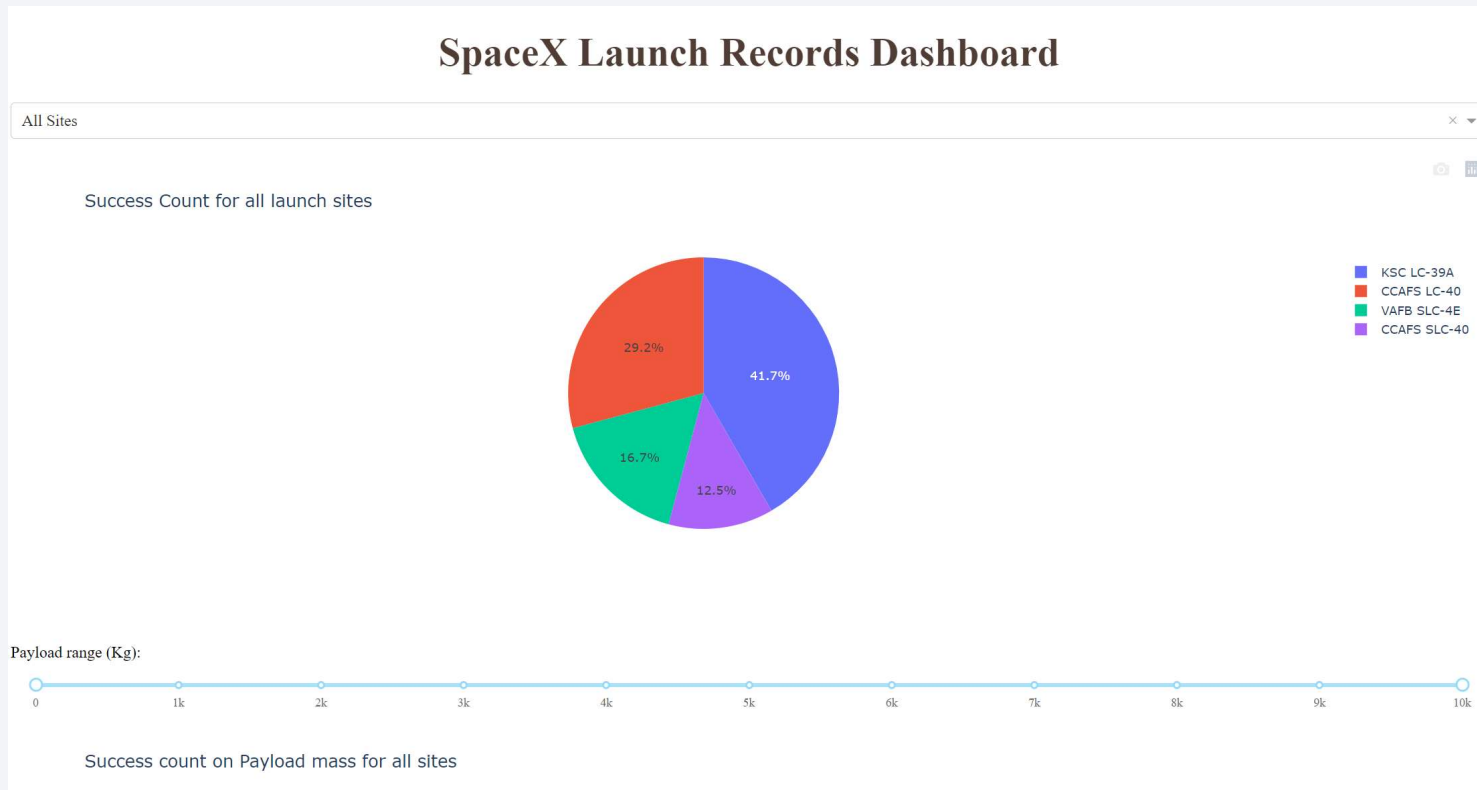# <Folium Map Screenshot 2>

# <Folium Map Screenshot 3>

Section 4
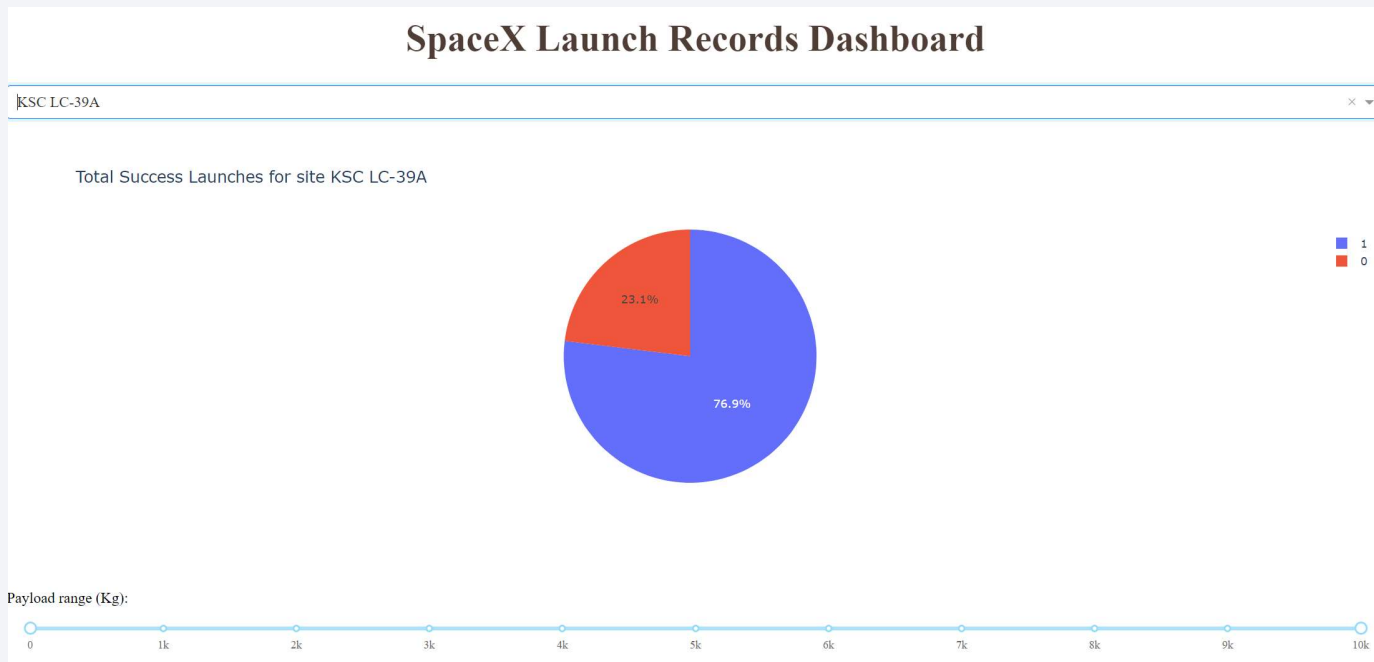
# Build a Dashboard with Plotly Dash

# <SpaceX Launch Records Dashboard>



KSC LC-39A amount to the highest portion compared to other sites

# <SpaceX Launch Records Dashboard Con>

**SpaceX Launch Records Dashboard**

KSC LC-39A

Total Success Launches for site KSC LC-39A

23.1%

76.9%

■ 1
■ 0

Payload range (Kg):

0    1k    2k    3k    4k    5k    6k    7k    8k    9k    10k

- KSC LC-39A have about 77% success rate.

# <SpaceX Launch Records Dashboard Con>



Success count on Payload mass for all sites

The best booster version is FT

Section 5

# Predictive Analysis (Classification)
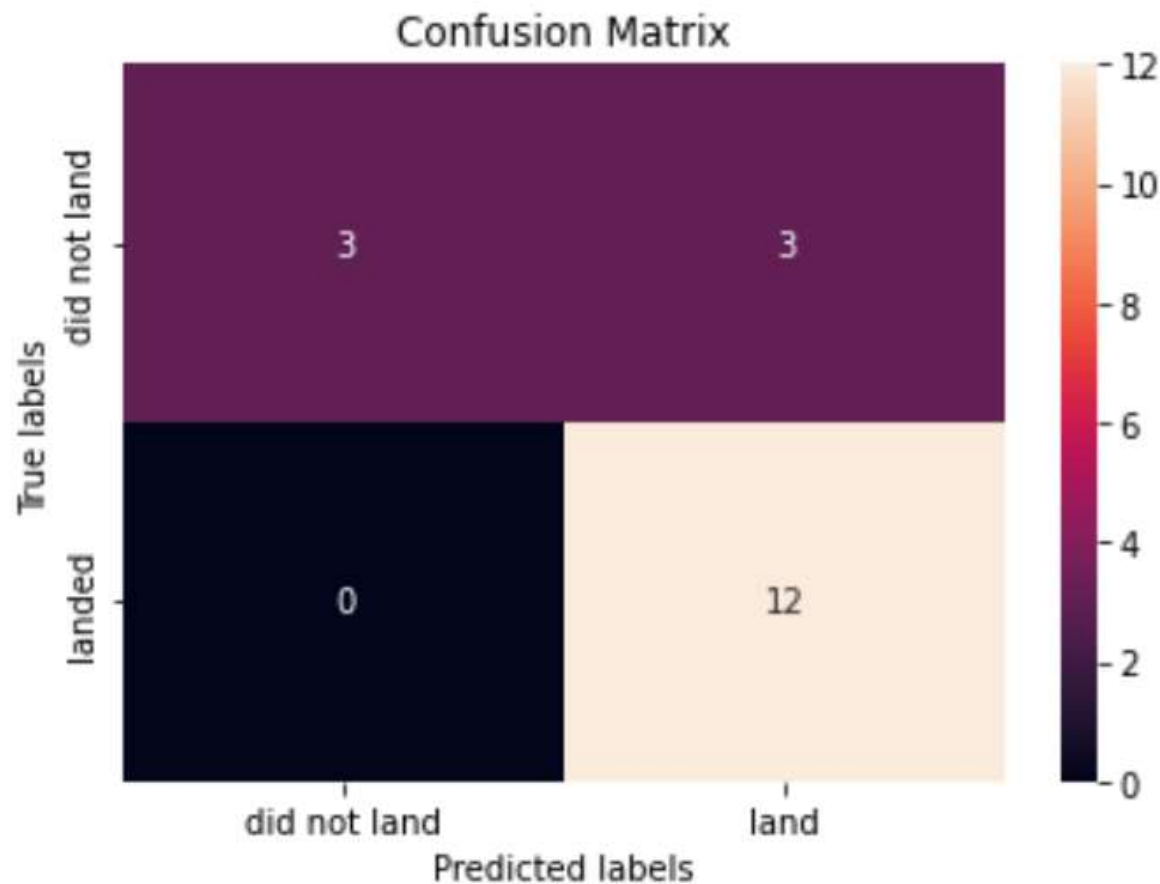
# Classification Accuracy

- For our case, all model presents equal accuracy, this may be due to split of training data set is random.

```
▶| list_score = [knn_cv.score(X_test,Y_test), tree_cv.score(X_test,Y_test), svm_cv.score(X_test,Y_test), logreg_cv.score(X_test,

    # Present the result
    df = pd.DataFrame(list_score, index=['KNN','Decision Tree','SVM','Logistic Regression'])
    df.columns = ['Score']

    df
```

0]:

|  | Score |
| --- | --- |
| KNN | 0.833333 |
| Decision Tree | 0.833333 |
| SVM | 0.833333 |
| Logistic Regression | 0.833333 |

# Confusion Matrix



Confusion Matrix

Examining the confusion matrix, we see that logistic regression can distinguish between the different classes. We see that the major problem is false positives.

# Conclusions

- Success rate of rocket launches increases over time

- Model selection is objective, and accuracy is very dependent on the data set.

- Introduce MLR in further study will tell us the weight of each factor in determining the success rate of each launch.

- False feature can be a problem.

# Thank you!