# Outline

- Executive Summary

- Introduction

- Methodology

- Results

- Conclusion

- Appendix

# Executive Summary

- Summary of methodologies

- Data Collection via API

- Data Collection via Web Scraping

- Data Wrangling

- Exploratory Data Analysis with Data Visualization

- Exploratory Data Analysis with SQL

- Interactive Visual Analytics with Folium

- Machine Learning Prediction


- Summary of all results

- Exploratory Data Analysis Results

- Interactive Analytics Demonstration in Screenshots

- Predictive Analysis Results

# 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!

- Which factors influence the landing outcome?

- What is the relationship between each variables and how does it affect the outcome?

- What is the best condition needed to increase the probability of successful landing?

Section 1

# Methodology

# Methodology

Executive Summary

- Data collection methodology:

    - Data was collected using SpaceX REST API and web scrapping from Wikipedia

- Perform data wrangling

    - 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

# Data Collection

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.

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

1. For REST API we used the get request. We decoded the response content as Json and turn it into a pandas dataframe using json_normalize(). We cleaned the data, checked for missing values and filled using the mean values.

2. For web scrapping, we used BeautifulSoup to extract the launch records as HTML table, parse the table and convert it into a pandas dataframe for further analysis.
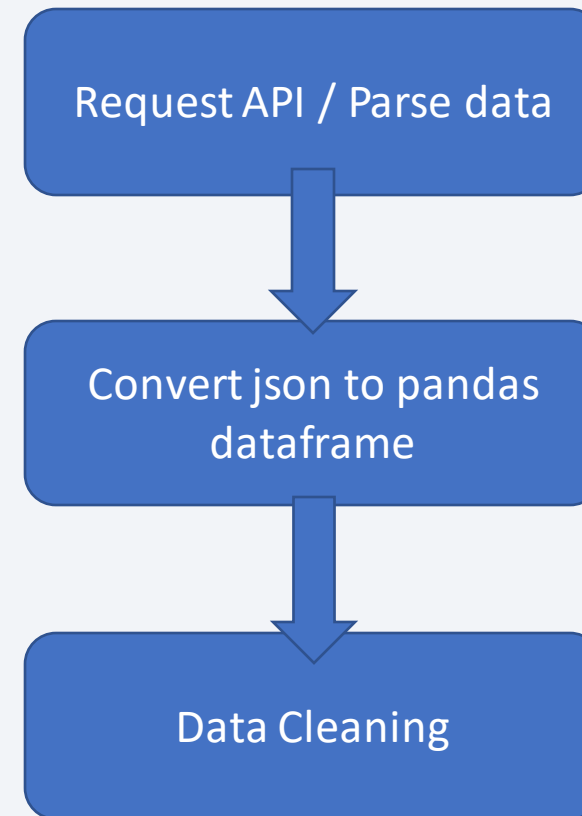
# Data Collection – SpaceX API

Step-by-step:

- 1 - Request API and parse the SpaceX launch data

- 2 - Use normalization method to convert json result into dataframe

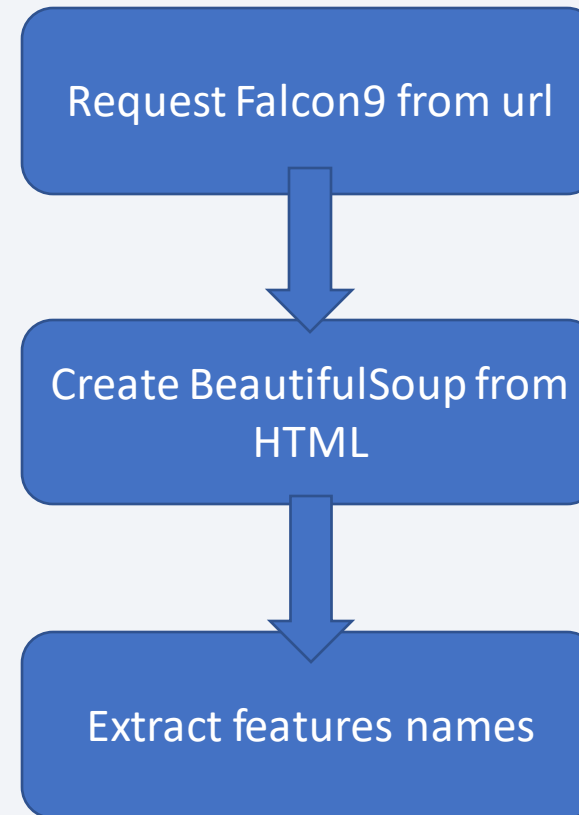- 3 - Performe data cleaning and fill missing value

Source Code GitHub



Request API / Parse data

↓

Convert json to pandas dataframe

↓

Data Cleaning

# Data Collection - Scraping

**Step-by-step:**

- 1 - Request the Falcon9 Launch Wiki page from url

- 2 - Create a BeautifulSoup from the HTML response

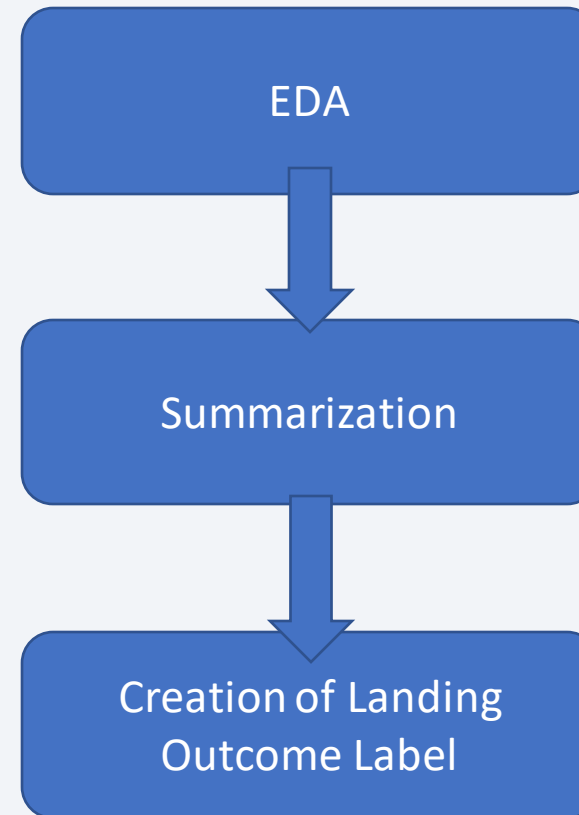- 3 - Extract all column/variable names from the HTML header

Source Code GitHub

```
Request Falcon9 from url
        ↓
Create BeautifulSoup from HTML
        ↓
Extract features names
```

# Data Wrangling

Step-by-step:

- 1 - Exploratory Data Analysis on the dataset

- 2 – Calculation of the number of launches on each site, then the number and occurrence of mission outcome per orbit type.

- 3 - Creation of a landing outcome label from the outcome column in order to facilitate further analysis and modelling.
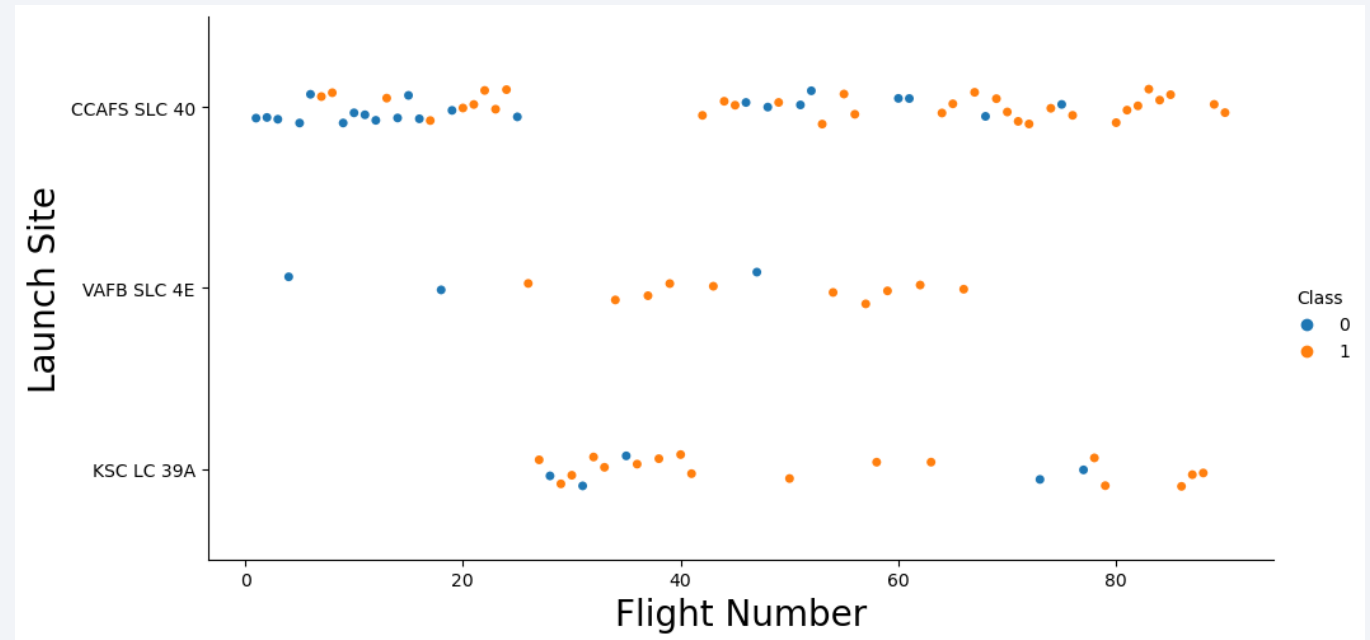
EDA

↓

Summarization

↓

Creation of Landing Outcome Label

Source Code GitHub

# EDA with Data Visualization

Scatter plots show dependency of attributes on each other. Once a pattern is determined from the graphs. It's very easy to see which factors affecting the most to the success of the landing outcomes.

**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

Source Code GitHub

# EDA with SQL

- With SQL we performed a series of queries to get insights of the dataset:

- - Displaying the names of the launch sites.

- - Displaying 5 records where launch sites begin with the string 'CCA'.

- - Displaying the total payload mass carried by booster launched by NASA (CRS).

- - Displaying the average payload mass carried by booster version F9 v1.1.

- - 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 sites

- names for in year 2015.

- - Rank the count of landing outcomes or success between the date 2010-06-04 and

- 2017-03-20, in descending order.

[Source Code GitHub](Source Code GitHub)

# Build an Interactive Map with Folium

- Markers for 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.

- Coloured Markers of the launch outcomes for each Launch Site:

- - Added coloured 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.
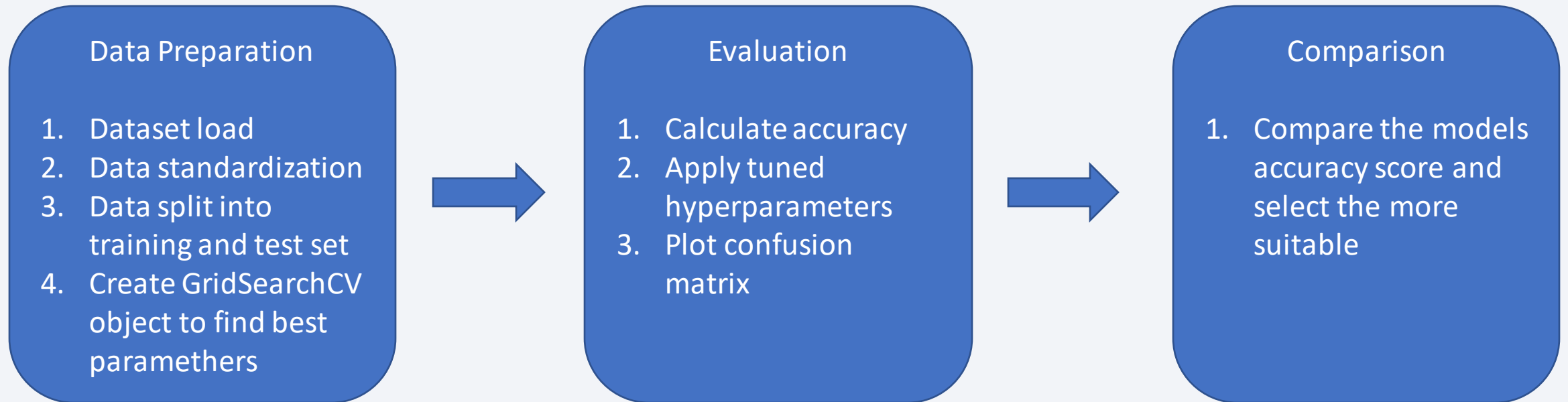
13

# 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.

Source Code GitHub

# Predictive Analysis (Classification)

**Data Preparation**

1. Dataset load
2. Data standardization
3. Data split into training and test set
4. Create GridSearchCV object to find best paramethers

**Evaluation**

1. Calculate accuracy
2. Apply tuned hyperparameters
3. Plot confusion matrix

**Comparison**

1. Compare the models accuracy score and select the more suitable

Four classification models were compared:

logistic regression, support vector machine, decision tree and K-NN.

15

Source Code GitHub

# Results

**Exploratory data analysis results:**

- SpaceX uses **4 launch sites**;

- The **first launches** were performed by **Space X** and **NASA**;

- The average payload of **F9 v1.1 booster is 2,928 kg**;

- The **first success landing outcome happened in 2015**;

- The **majority 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;
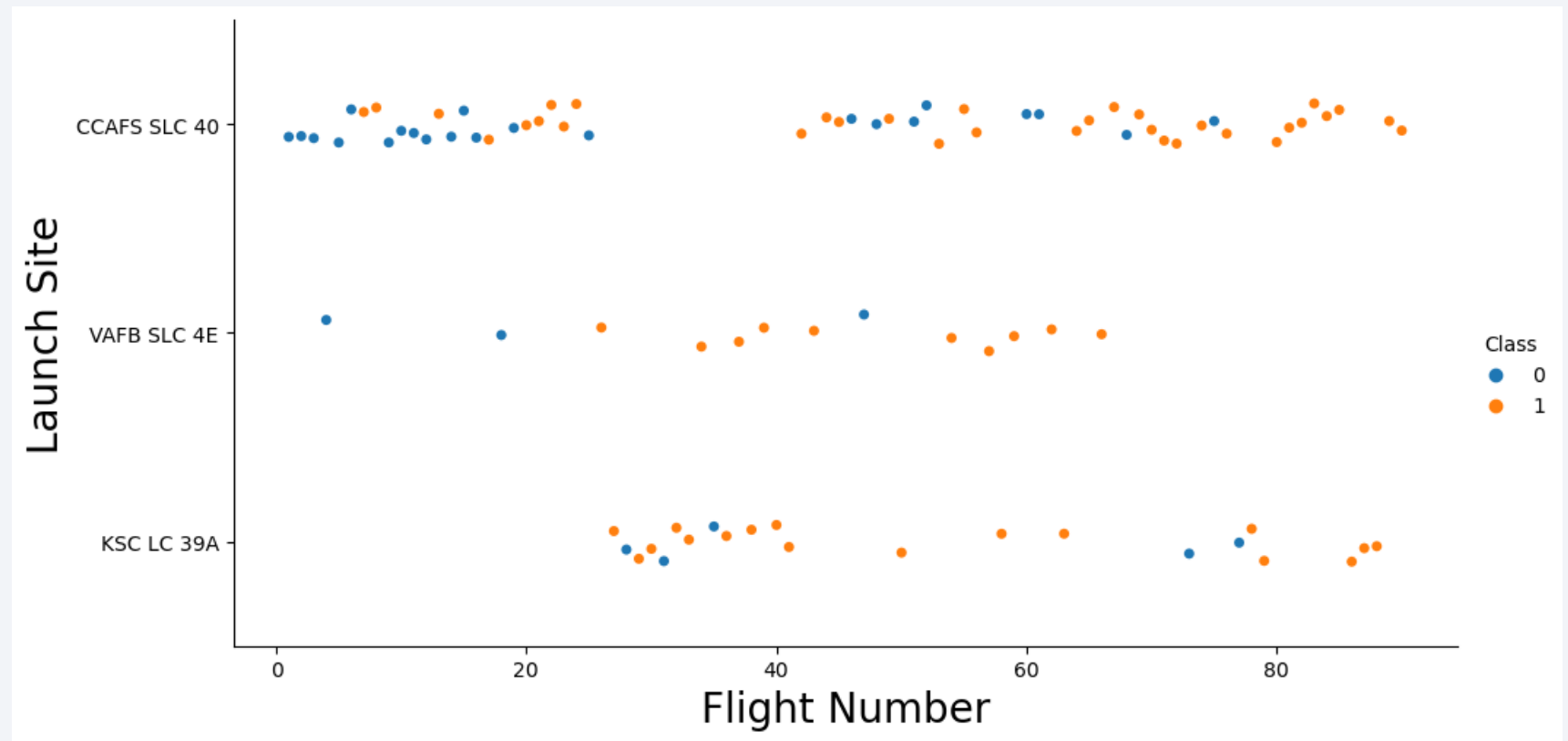
Source Code GitHub

Section 2

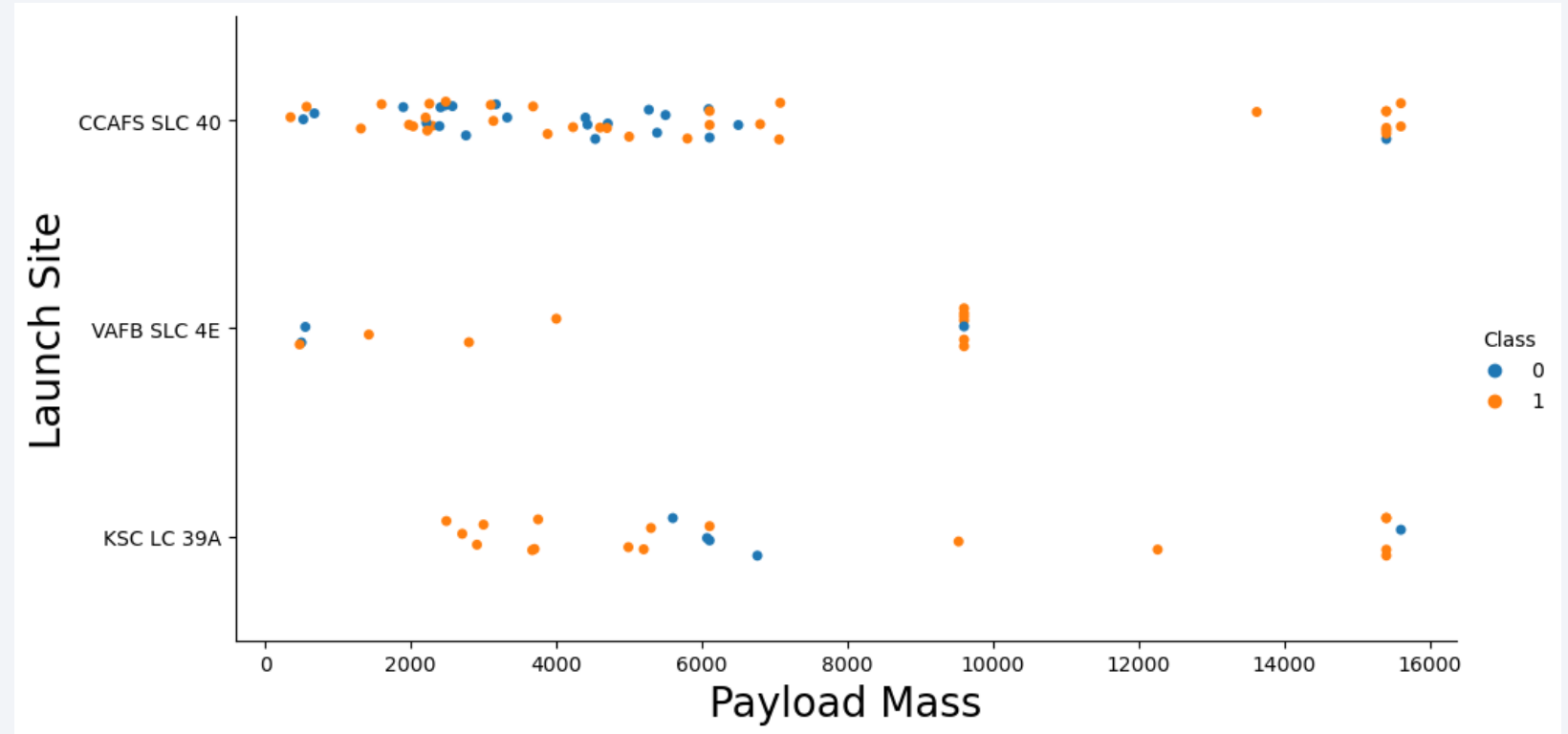# Insights drawn from EDA

# 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.

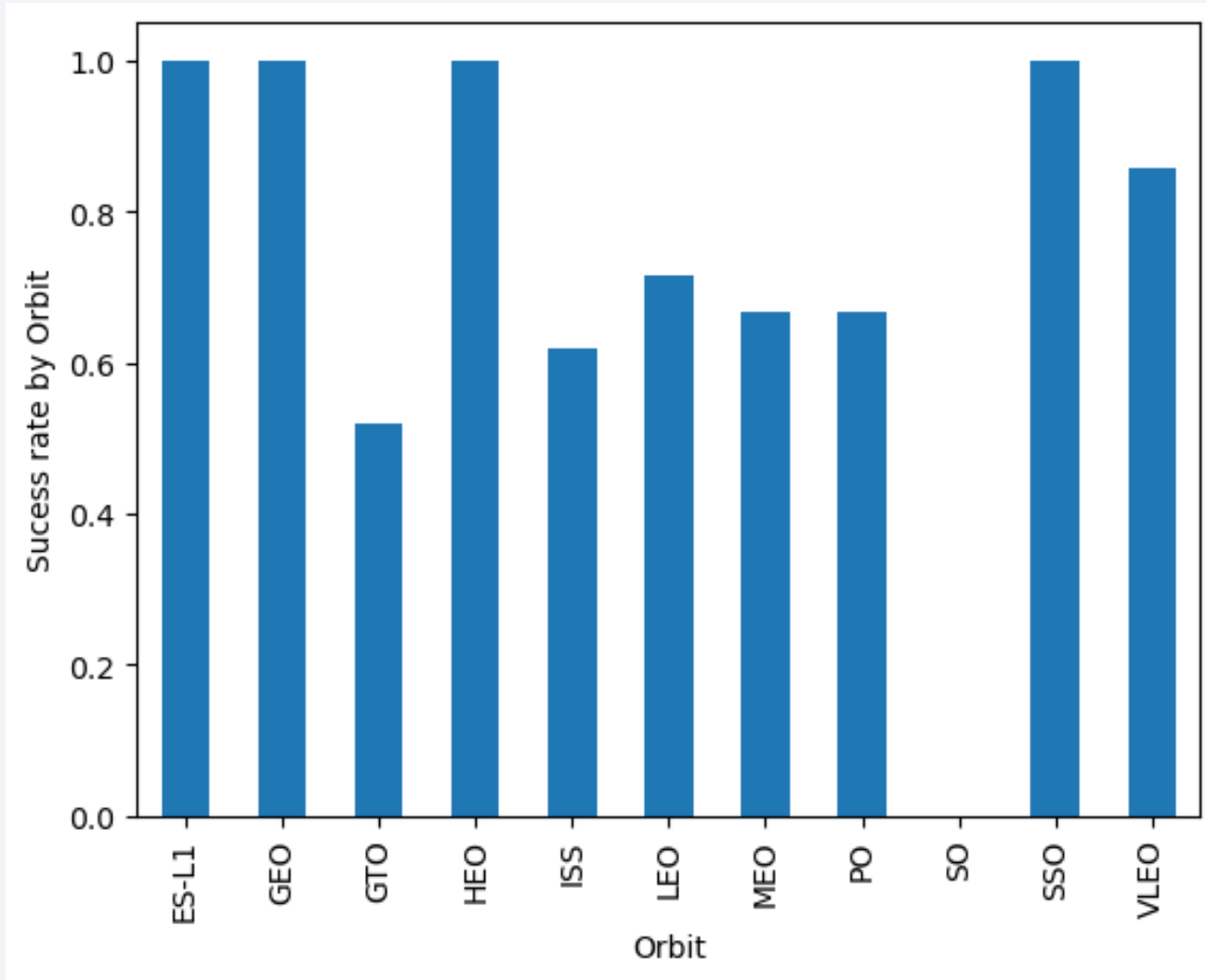- It can be assumed that each new launch has a higher rate of success.

# 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.

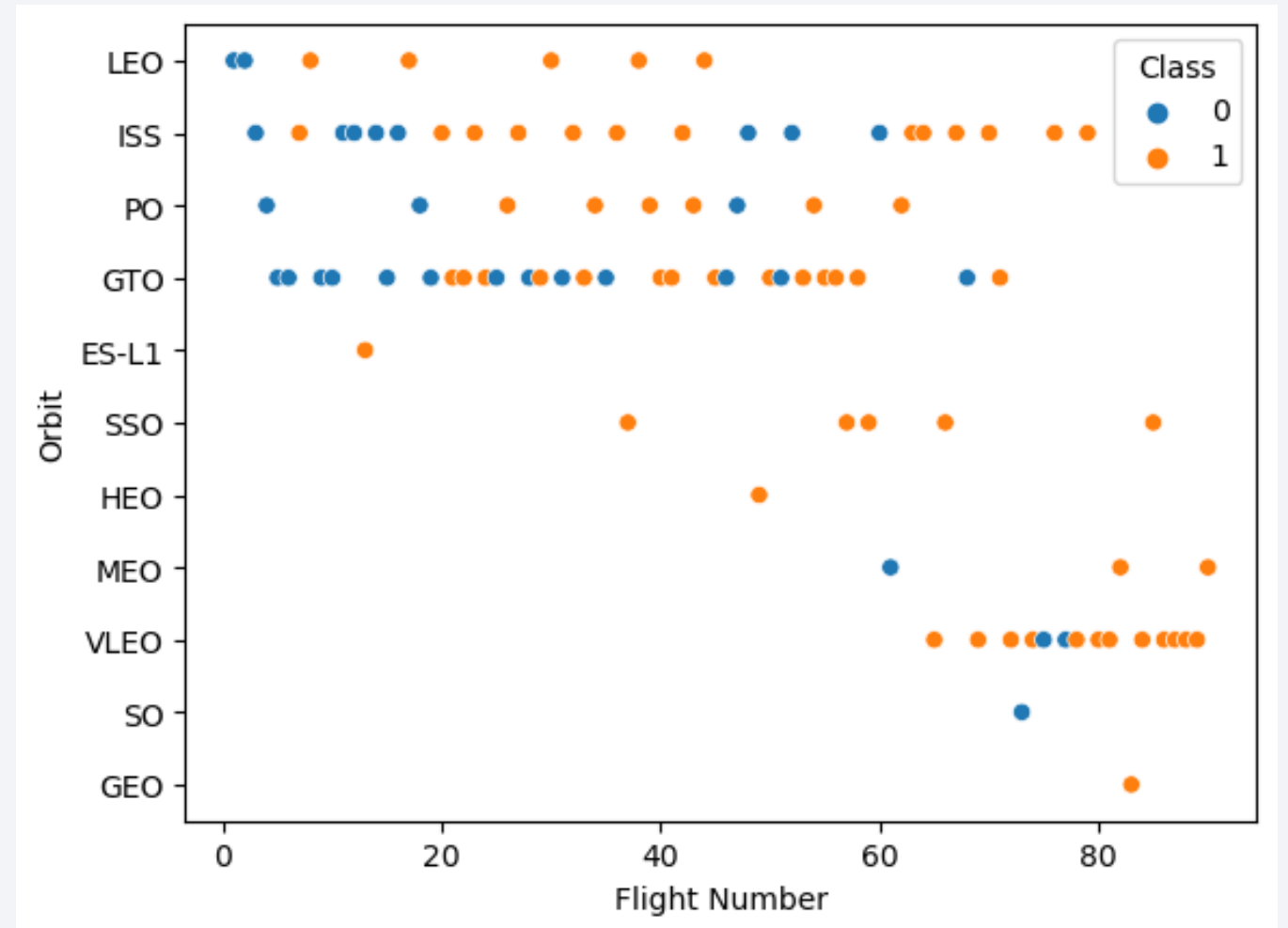- KSC LC 39A has a 100% success rate for payload mass under 5500 kg too.

# 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%:
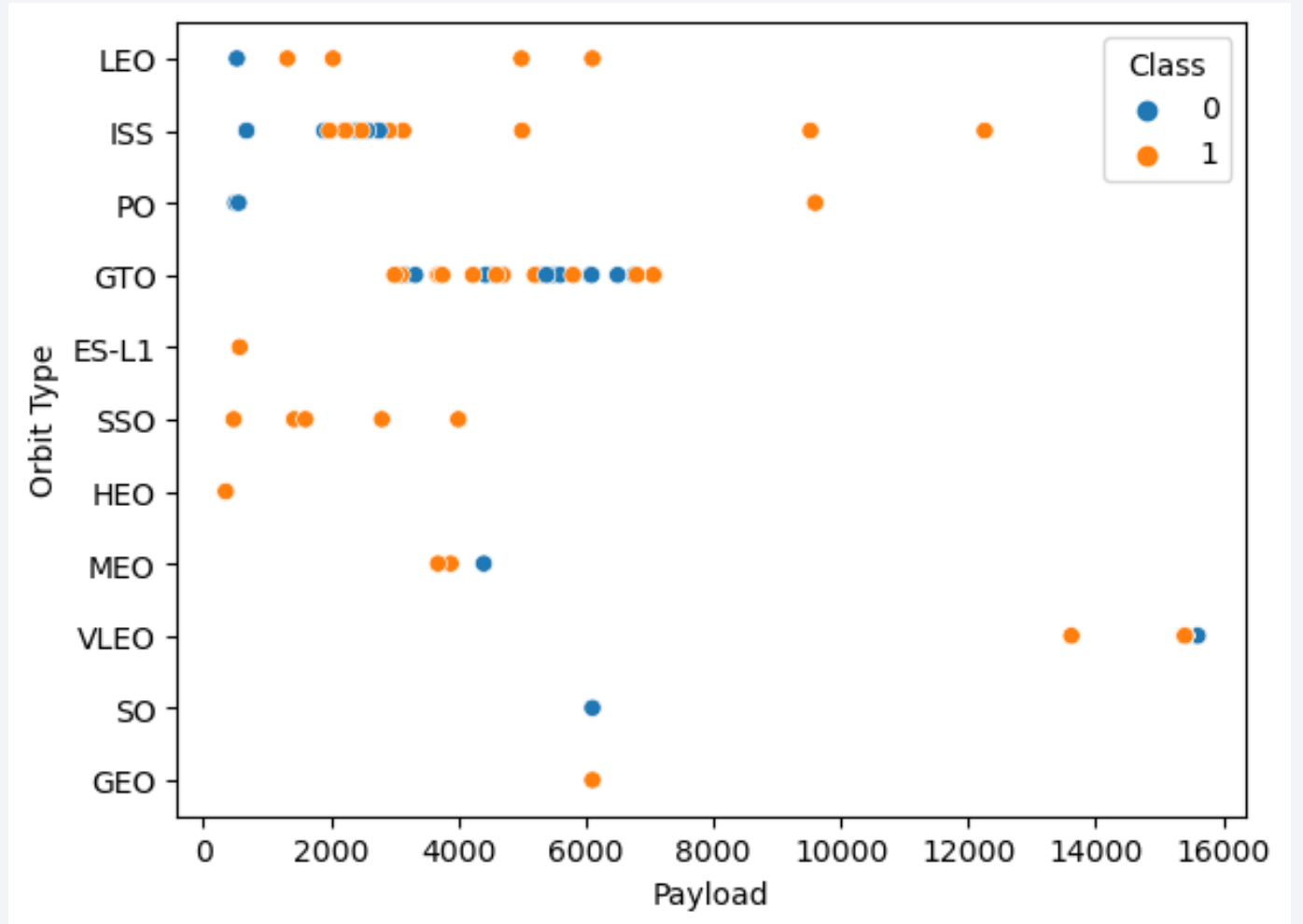
  - GTO, ISS, LEO, MEO, PO

# 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.
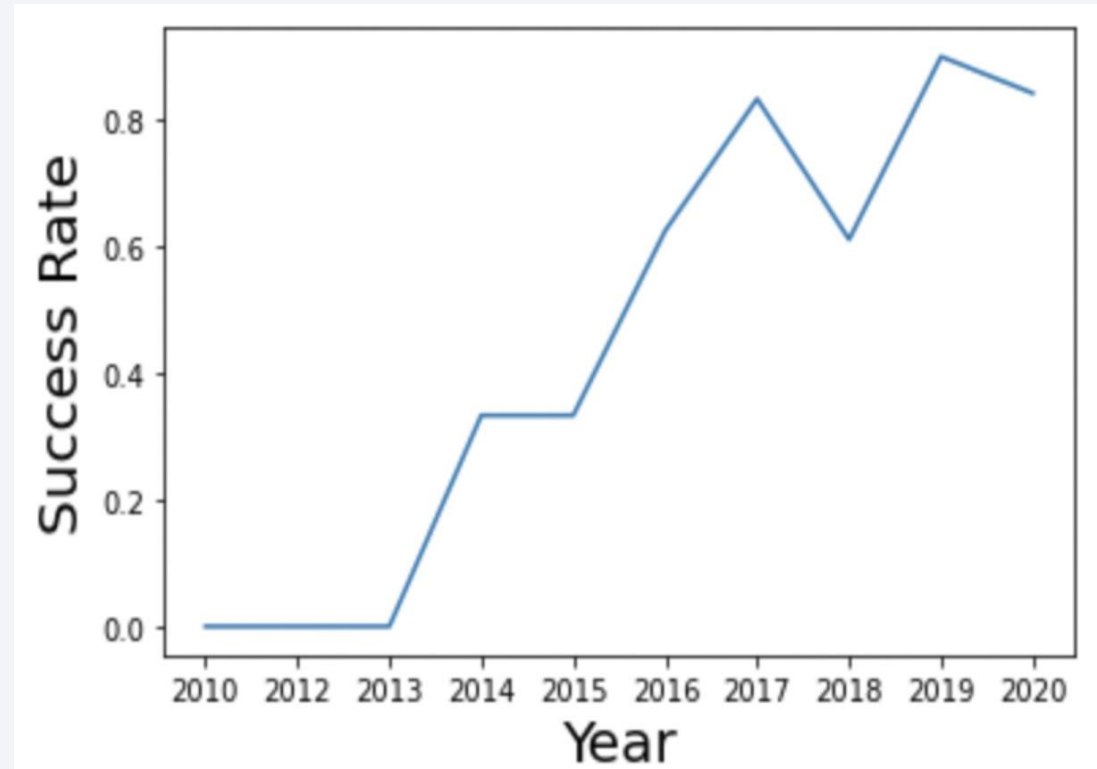
# Payload vs. Orbit Type

- 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

We used the key word DISTINCT to show only unique
launch sites from the SpaceX data.

```
In [4]: %sql select distinct launch_site from SPACEXDATASET;

         * ibm_db_sa://wzf08322:***@0c77d6f2-5da9-48a9-81f8-86b520b87518.bs2io90l08kqb1od8lcg.databases.appdomain.cloud:31198/bludb
        Done.

Out[4]:
```

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

# Launch Site Names Begin with 'CCA'

We used the query above to display 5 records where launch sites begin with CCA

In [5]: `%sql select * from SPACEXDATASET where launch_site like 'CCA%' limit 5;`

* ibm_db_sa://wzf08322:***@0c77d6f2-5da9-48a9-81f8-86b520b87518.bs2io90l08kqb1od8lcg.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 |

# Total Payload Mass

We calculated the total payload carried by boosters
from NASA as 45596 using the query below

```
In [6]: %sql select sum(payload_mass__kg_) as total_payload_mass from SPACEXDATASET where customer = 'NASA (CRS)';
         * ibm_db_sa://wzf08322:***@0c77d6f2-5da9-48a9-81f8-86b520b87518.bs2io90l08kqb1od8lcg.databases.appdomain.cloud:31198/bludb
        Done.
Out[6]:
        total_payload_mass
        45596
```

# Average Payload Mass by F9 v1.1

We calculated the average payload mass carried by
booster version F9 v1.1 as 2928.4

```
In [7]:  %sql select avg(payload_mass__kg_) as average_payload_mass from SPACEXDATASET where booster_version like '%F9 v1.1%';

         * ibm_db_sa://wzf08322:***@0c77d6f2-5da9-48a9-81f8-86b520b87518.bs2io90108kqb1od8lcg.databases.appdomain.cloud:31198/bludb
         Done.
```

Out[7]:

| average_payload_mass |
|----------------------|
| 2534                 |

# First Successful Ground Landing Date

We use the min() function to find the result
We observed that the dates of the first successful landing
outcome on ground pad was 22nd December 2015

```
In [8]: %sql select min(date) as first_successful_landing from SPACEXDATASET where landing__outcome = 'Success (ground pad)';

         * ibm_db_sa://wzf08322:***@0c77d6f2-5da9-48a9-81f8-86b520b87518.bs2io90l08kqb1od8lcg.databases.appdomain.cloud:31198/bludb
        Done.

Out[8]:   first_successful_landing

          2015-12-22
```

# Successful Drone Ship Landing with Payload between 4000 and 6000

We used the WHERE clause to filter for boosters which have successfully landed on drone ship and applied the AND condition to determine successful landing with payload mass greater than 4000 but less than 6000

```
In [9]: %sql select booster_version from SPACEXDATASET where landing__outcome = 'Success (drone ship)' and payload_mass__kg_ between 4000 and 6000;

 * ibm_db_sa://wzf08322:***@0c77d6f2-5da9-48a9-81f8-86b520b87518.bs2io90l08kqb1od8lcg.databases.appdomain.cloud:31198/bludb
Done.
```

Out[9]:

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

# Total Number of Successful and Failure Mission Outcomes

We used like % to filter for WHERE MissionOutcome
was a success or a failure.

```
In [10]: %sql select mission_outcome, count(*) as total_number from SPACEXDATASET group by mission_outcome;

          * ibm_db_sa://wzf08322:***@0c77d6f2-5da9-48a9-81f8-86b520b87518.bs2io90108kqb1od8lcg.databases.appdomain.cloud:31198/bludb
         Done.
```

Out[10]:

| mission_outcome | total_number |
|---|---|
| Failure (in flight) | 1 |
| Success | 99 |
| Success (payload status unclear) | 1 |

# Boosters Carried Maximum Payload

We determined the booster that have carried the maximum payload using a subquery in the WHERE clause and the MAX() function.

```
In [11]: %sql select booster_version from SPACEXDATASET where payload_mass__kg_ = (select max(payload_mass__kg_) from SPACEXDATASET);
         * ibm_db_sa://wzf08322:***@0c77d6f2-5da9-48a9-81f8-86b520b87518.bs2io90108kqb1od8lcg.databases.appdomain.cloud:31198/bludb
         Done.
Out[11]:
```

| 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

We used a combinations of the WHERE clause, LIKE,
AND, and BETWEEN conditions to filter for failed landing
outcomes in drone ship, their booster versions, and
launch site names for year 2015

```
In [12]:  %%sql select monthname(date) as month, date, booster_version, launch_site, landing__outcome from SPACEXDATASET
          where landing__outcome = 'Failure (drone ship)' and year(date)=2015;

 * ibm_db_sa://wzf08322:***@0c77d6f2-5da9-48a9-81f8-86b520b87518.bs2io90108kqb1od8lcg.databases.appdomain.cloud:31198/bludb
Done.
```

Out[12]:

| MONTH | DATE | booster_version | launch_site | landing__outcome |
|---|---|---|---|---|
| January | 2015-01-10 | F9 v1.1 B1012 | CCAFS LC-40 | Failure (drone ship) |
| April | 2015-04-14 | F9 v1.1 B1015 | CCAFS LC-40 | Failure (drone ship) |

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

We selected Landing outcomes and the COUNT of landing outcomes from the data and used the WHERE clause to filter for landing outcomes BETWEEN 2010-06-04 to 2010-03-20.

We applied the GROUP BY clause to group the landing outcomes and the ORDER BY clause to order the grouped landing outcome in descending order.

```
In [13]:  %%sql select landing__outcome, count(*) as count_outcomes from SPACEXDATASET
          where date between '2010-06-04' and '2017-03-20'
          group by landing__outcome
          order by count_outcomes desc;
```

 * ibm_db_sa://wzf08322:***@0c77d6f2-5da9-48a9-81f8-86b520b87518.bs2io90108kqb1od8lcg.databases.appdomain.cloud:31198/bludb
Done.

Out[13]:

| landing__outcome | count_outcomes |
|---|---|
| No attempt | 10 |
| Failure (drone ship) | 5 |
| Success (drone ship) | 5 |
| Controlled (ocean) | 3 |
| Success (ground pad) | 3 |
| Failure (parachute) | 2 |
| Uncontrolled (ocean) | 2 |
| Precluded (drone ship) | 1 |

Section 3

# Launch Sites Proximities Analysis

# All launch sites



VAFB
Los Angeles
SLC-
4E

CCAFS
SLC-
40A
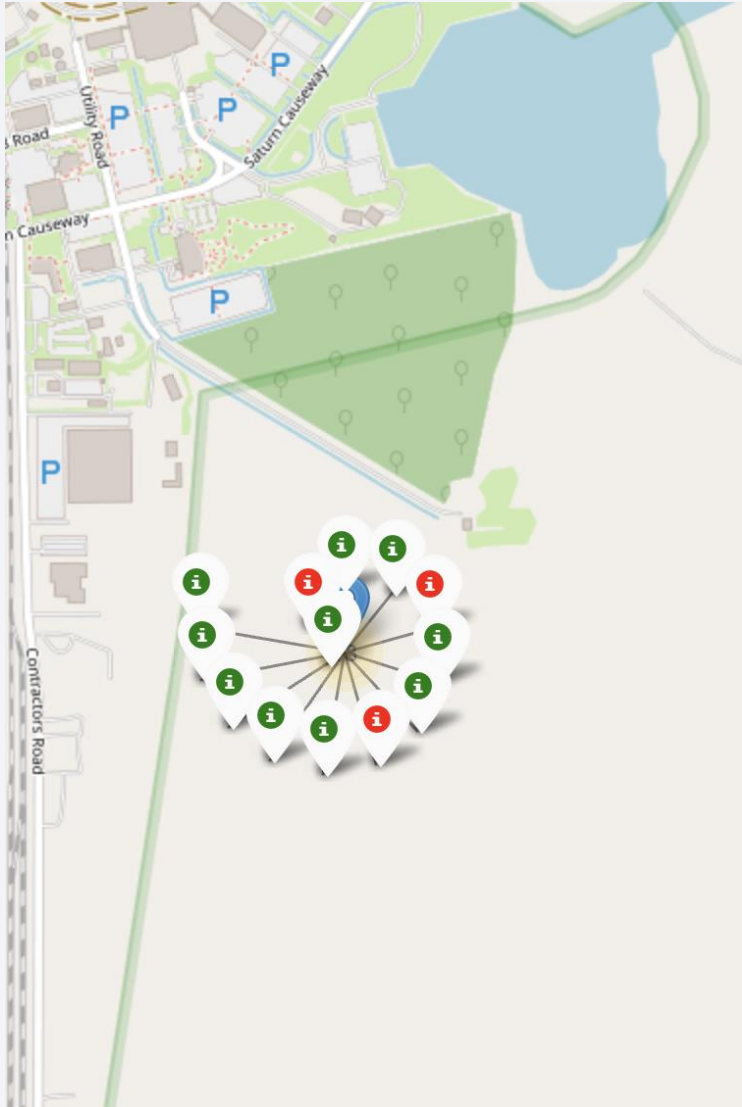
Launch sites are located
by two different oceans,
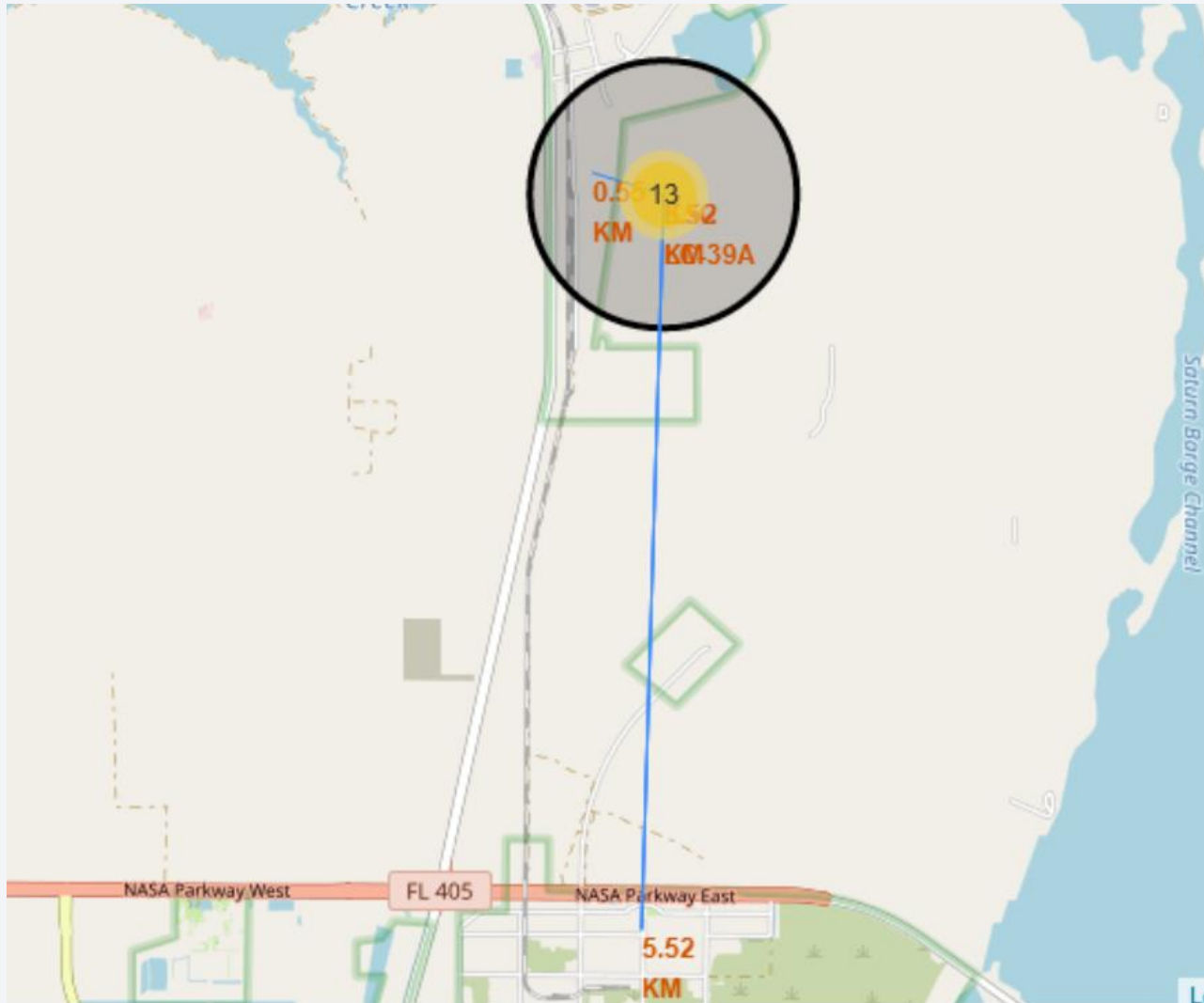probably for security reasons 35

# Color labeled launch sites

Green Marker = Successful Launch
Red Marker = Failed Launch

Launch Site KSC LC-39A has a very high Success Rate.

# KSC LC-39A Logistics Aspects



Launch site KSC LC-39A has good logistics aspects, being near railroad and road and relatively far from inhabited areas.

Section 4

# Build a Dashboard
# with Plotly Dash
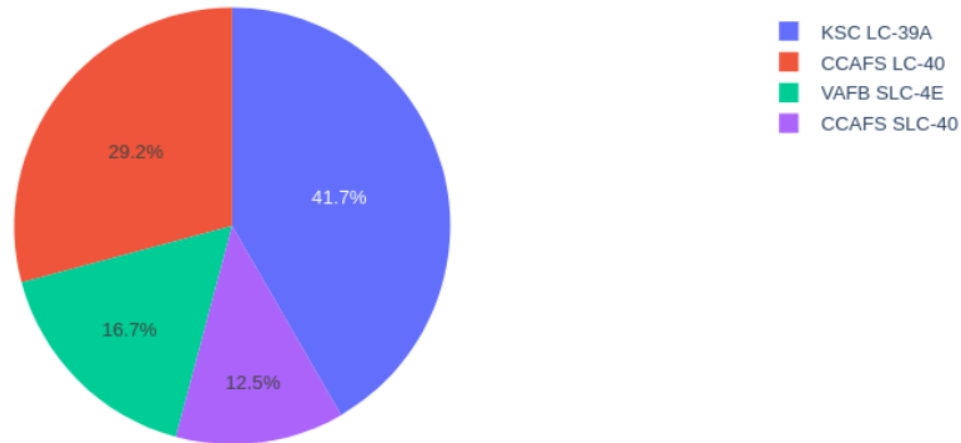
# Success percentage by each site
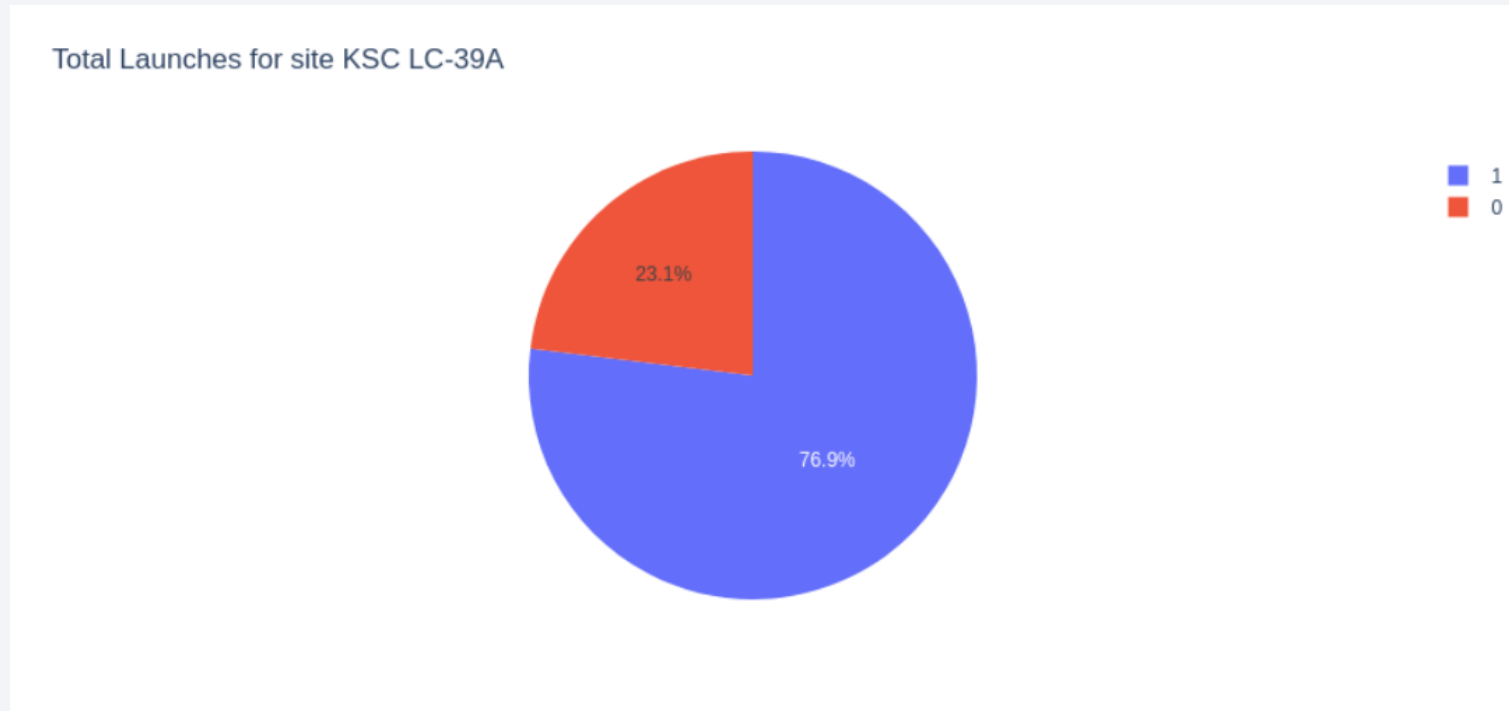


The chart shows that from all the sites, KSC LC-39A (41,7%) has the most successful launches, followed by CCAFS LC-40.

# The highest launch success ratio
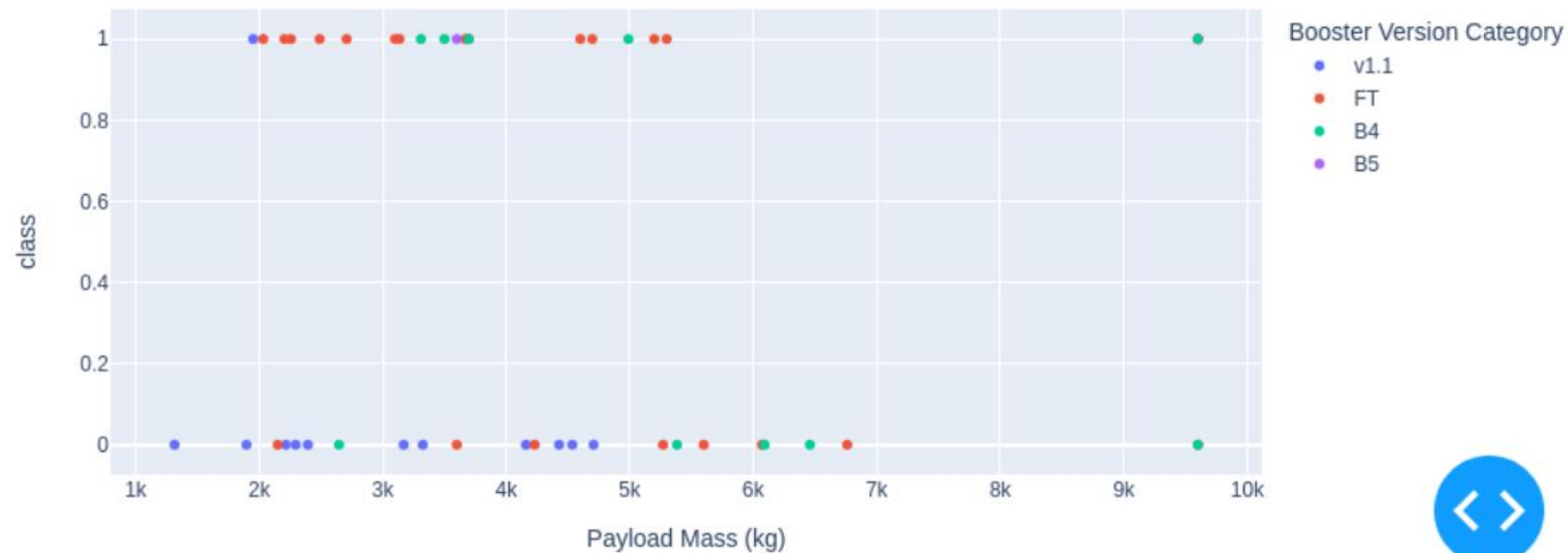


76.9% of launches are successful at KSC LC-39A

# Payload vs Launch



Payloads under 6,000kg and FT boosters are the most successful combination.

Section 5

# Predictive Analysis (Classification)

# Classification Accuracy

```python
print("Model\t\tAccuracy\tTestAccuracy")#,logreg_cv.best_score_)
print("LogReg\t\t{}\t\t{}".format((logreg_cv.best_score_).round(5), logreg_cv.score(X_test, Y_test).round(5)))
print("SVM\t\t{}\t\t{}".format((svm_cv.best_score_).round(5), svm_cv.score(X_test, Y_test).round(5)))
print("Tree\t\t{}\t\t{}".format((tree_cv.best_score_).round(5), tree_cv.score(X_test, Y_test).round(5)))
print("KNN\t\t{}\t\t{}".format((knn_cv.best_score_).round(5), knn_cv.score(X_test, Y_test).round(5)))

comparison = {}

comparison['LogReg'] = {'Accuracy': logreg_cv.best_score_.round(5), 'TestAccuracy': logreg_cv.score(X_test, Y_test).round(5)}
comparison['SVM'] = {'Accuracy': svm_cv.best_score_.round(5), 'TestAccuracy': svm_cv.score(X_test, Y_test).round(5)}
comparison['Tree'] = {'Accuracy': tree_cv.best_score_.round(5), 'TestAccuracy': tree_cv.score(X_test, Y_test).round(5)}
comparison['KNN'] = {'Accuracy': knn_cv.best_score_.round(5), 'TestAccuracy': knn_cv.score(X_test, Y_test).round(5)}
```

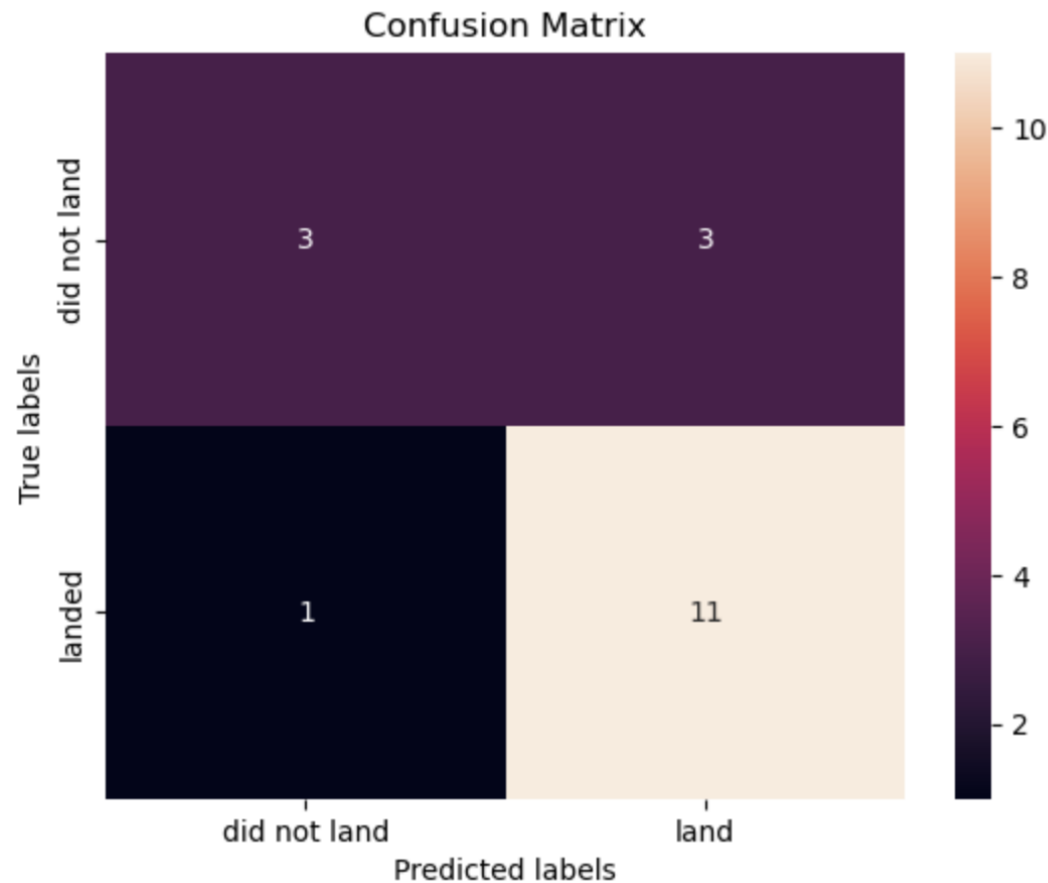| Model  | Accuracy | TestAccuracy |
|--------|----------|--------------|
| LogReg | 0.84643  | 0.83333      |
| SVM    | 0.84821  | 0.83333      |
| Tree   | 0.88929  | 0.77778      |
| KNN    | 0.84821  | 0.83333      |

The model with the highest classification accuracy is Decision
Tree Classifier, which displays an accuracy over than 88%

# Confusion Matrix

```
In [34]:  yhat = tree_cv.predict(X_test)
          plot_confusion_matrix(Y_test,yhat)
```



Examining the confusion matrix, we see that Decision Tree can distinguish between the different classes. We see that the major problem is false positives (3).

# 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.
- KSC LC-39A has the highest success rate of the launches
- from all the sites.
- Orbits ES-L1, GEO, HEO and SSO have 100% success rate.

Thank you!