



SpaceX Data Exploration And Predictions

*Ejale Emmanuel Osobase
September 2022*

OUTLINE

- ▶ Executive Summary
- ▶ Introduction
- ▶ Methodology
- ▶ Results
 - ▶ Visualization – Charts
 - ▶ Dashboard
- ▶ Discussion
 - ▶ Findings & Implications
- ▶ Conclusion
- ▶ Appendix



EXECUTIVE SUMMARY

- ▶ Summary of Methodology
 - ▶ Data collection with API and Web Scraping
 - ▶ Data Wrangling
 - ▶ Exploratory Data Analysis with SQL
 - ▶ Exploratory Data Analysis with Data Visualization
 - ▶ Interactive Visual Analysis with Folium
 - ▶ Interactive Visual Analysis with Plotly Dash
- ▶ Predictive Analysis
 - ▶ Machine Learning
- ▶ Result Summary
 - ▶ Exploratory Data Analysis
 - ▶ Visuals
 - ▶ Predictive Analysis

INTRODUCTION

- Project Scenario Overview

SpaceX Falcon 9 launches with a cost of 62 million dollars while other providers cost upwards of 165 million dollars each, this is due to SpaceX being able to reuse the first stage of the rocket launched. Therefore, if we can predict successful landing of the first stage we can determine the cost of a launch.

- Questions

- Determine if SpaceX will reuse the first stage
- Determine the price of each Launch
- Train ML models and use public information to predict if SpaceX will reuse the first stage



METHODOLOGY

- ▶ Data Collection:
 - ▶ Using SpaceX API
 - ▶ Web Scraping from Wikipedia with BeautifulSoup
- ▶ Data Wrangling
 - ▶ Exploratory Data Analysis (EDA) with SQL and Visuals
 - ▶ Interactive visual analysis with folium and Plotly Dash
- ▶ Predictive Analysis:
 - ▶ Train classification models for predictions
 - ▶ Using folium world map to calculate proximity to city, coastline etc.

Data Collection- SpaceX API

- ▶ The data was collected with get request from SpaceX API
- ▶ Use the API to get information about the launches
- ▶ Apply accurate names to each columns
- ▶ Combine the columns into a dictionary
- ▶ Create a pandas data frame from the dictionary
- ▶ Filter the data frame to only include falcon 9 launches

Data Collection- Web Scraping

- ▶ Request the Falcon 9 launch Wikipedia page from its URL
- ▶ Extract Falcon 9 records HTML tables from Wikipedia with BeautifulSoup
- ▶ Parse the table and convert it into a Pandas Data frame

Data Wrangling

- ▶ Remove null values
- ▶ Identify numerical and categorical columns
- ▶ Calculate the number of launches on each site
- ▶ Calculate the number of occurrences of each orbit
- ▶ Calculate occurrences of mission outcome per orbit type
- ▶ Create a landing Outcome column

EDA with SQL Analysis

- ▶ Space X dataset is Loaded and formatted within IBM Watson DB2 database and explored using jupyter notebook.
 - ▶ Exploration carried out using SQL to derive insights. These queries include:
 1. Names of unique Launch sites in the space mission
 2. Display 5 records where launch sites begin with the string 'CCA'
 3. Display the total payload mass carried by booster launched by NASA (CRS)
 4. Display average payload mass carried by booster version F9 v1.1
 5. List the date when the first successful landing outcome in ground pad was achieved

Link to the Notebook: https://github.com/EmmaOsobase/coursera-Capstone-Project/blob/main/jupyter-labs-eda-sql-coursera_sqlite.ipynb

EDA And Interactive Visual Analytics

- ▶ Display Unique Launch sites
- ▶ Total Mass carried by NASA (CRS)
- ▶ Average Payload carried by F9 v1.1
- ▶ First successful ground pad landing date
- ▶ Names of Booster versions which have carried maximum payload mass.
- ▶ How Flight Number and Payload affects the landing outcome
- ▶ Relationship between Flight Number and Launch site
- ▶ Relationship between Payload and Launch site
- ▶ How the Orbit type affects the success rate
- ▶ Yearly landing success rate

Plotly Dash dashboard Analysis

- ▶ Created an interactive webpage dashboard using plotly dash
- ▶ Plot a pie chart for success launches for selected launch sites
- ▶ Scattered chart for success count to show how payload mass and class for selected sites, and colour-label points for booster-versions. This chart is to show how payload mass affects the mission outcome for the selected site.

Interactive map with Folium

- ▶ Launch Sites Location Analysis using Folium
- ▶ Marked all Launch sites on the map using NASA Johnson Space center coordinates as the start location
- ▶ Add Circles and Markers to the map to pin point and label launch site using launch sites coordinates.
- ▶ Marked successful/failed launches for each site on the map
- ▶ Use marker cluster to highlight successful and failed launches based on colour.
- ▶ Calculate launch site distance from various proximities; such as coastline, Railway tracks, city and Highway.

Predictive Analysis

- ▶ Import the necessary Libraries
- ▶ Create a function to plot the confusion matrix
- ▶ Create numpy array from the column class and assign as my dependent variable
- ▶ Standardize fit and transform the independent variable
- ▶ Split the data to train and test the datasets
- ▶ Create classification objects, using LR, SVM, and decision tree, then create a GridSearch and fit to find the best parameter
- ▶ Calculate objects accuracy on the test dataset using method score



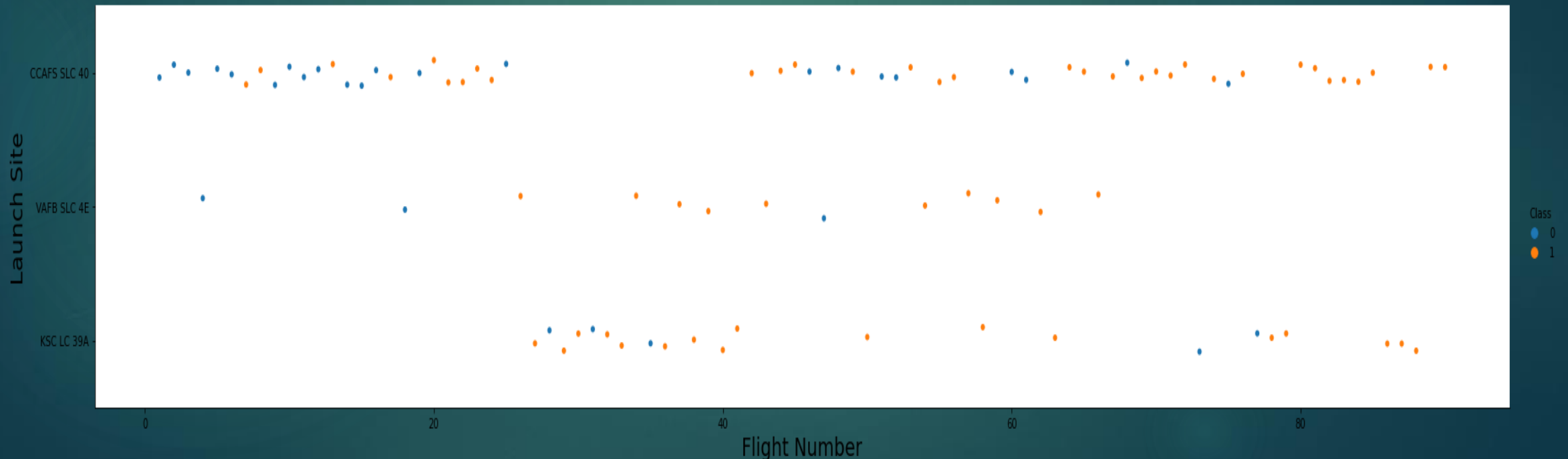
RESULTS

- ▶ EDA with visualization
- ▶ EDA with SQL
- ▶ Interactive map with Folium
- ▶ Plotly Dash dashboard
- ▶ Predictive analysis

EDA with Visualization

Flight Number Vs Launch Site

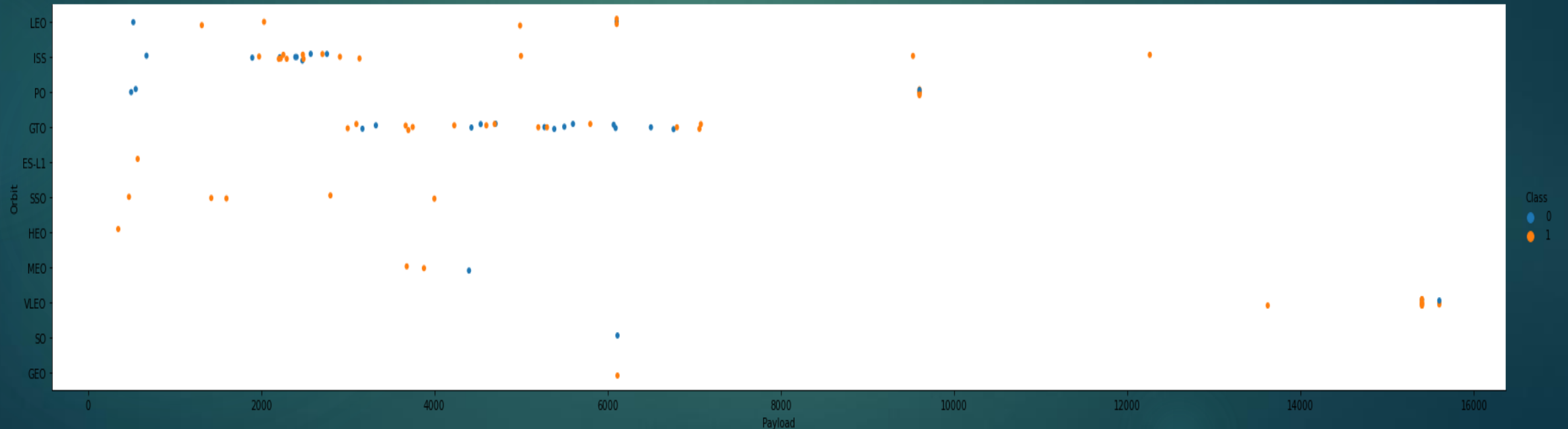
From this scattered plot we can see that irrespective of the launch site, the more flights are made the higher the success rate



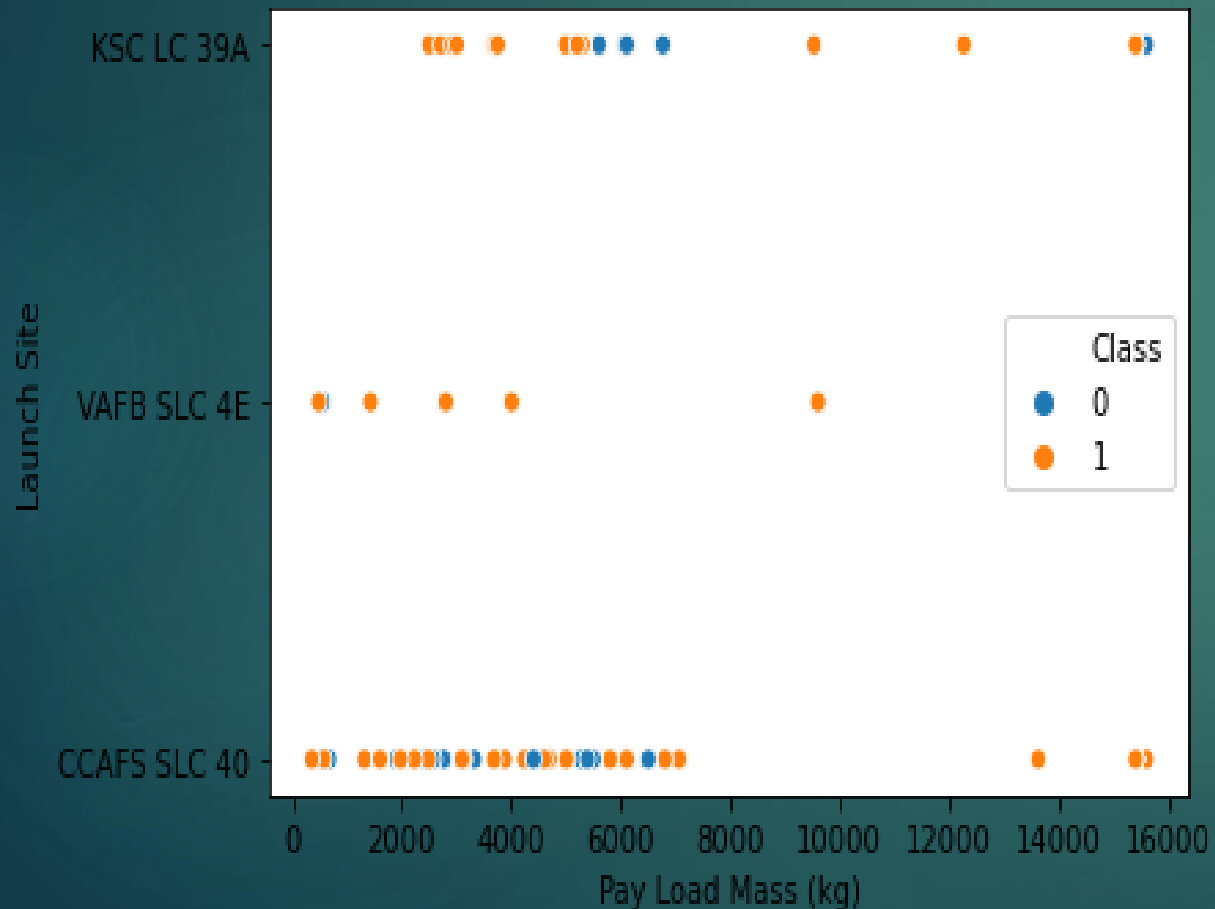
EDA with Visualization

Payload Vs Orbit

From this scattered plot we can see that with heavier payload the successful landing rate are more for Polar, LEO and ISS. However GTO has no distinguishable difference as both successful and failed landing outcomes are there.



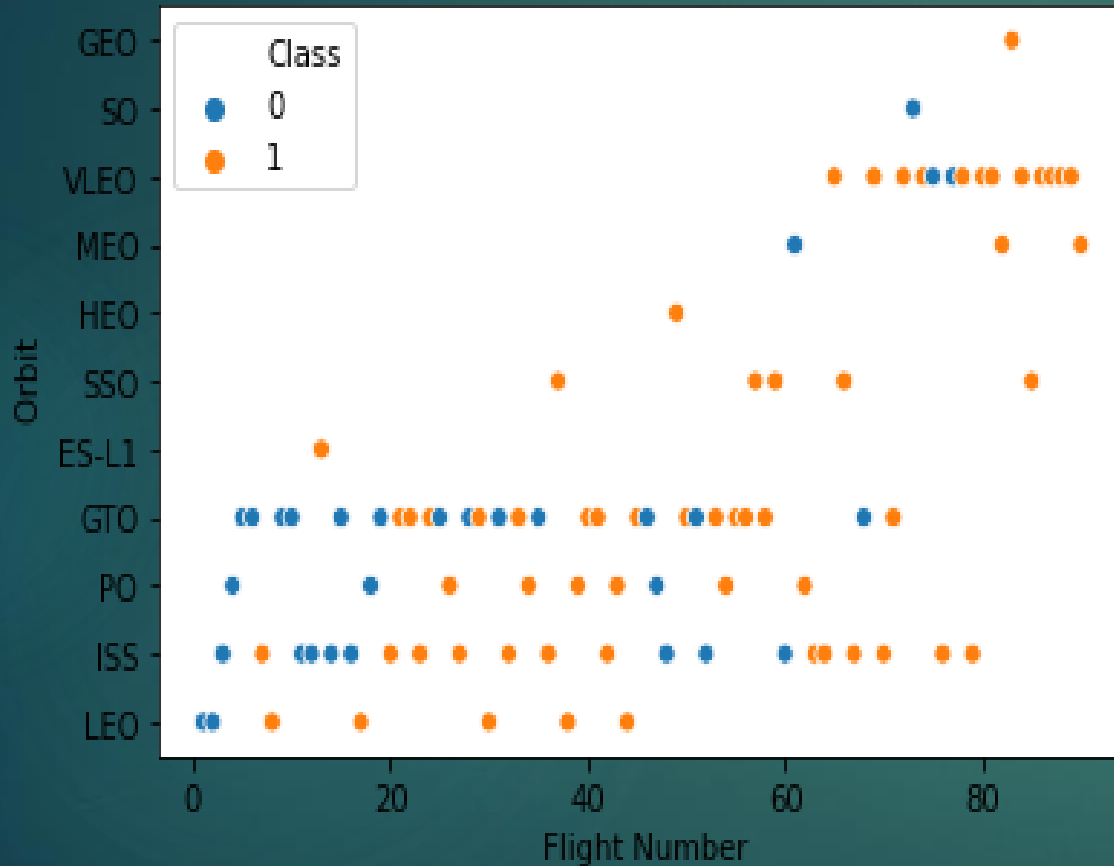
EDA with Visualization



Payload Vs Launch Site

In this scattered chart of Payload Vs launch site we can see that irrespective of the launch sites, rockets with payload $> 4000\text{kg}$ and $< 6000\text{kg}$ have a higher failure rate.

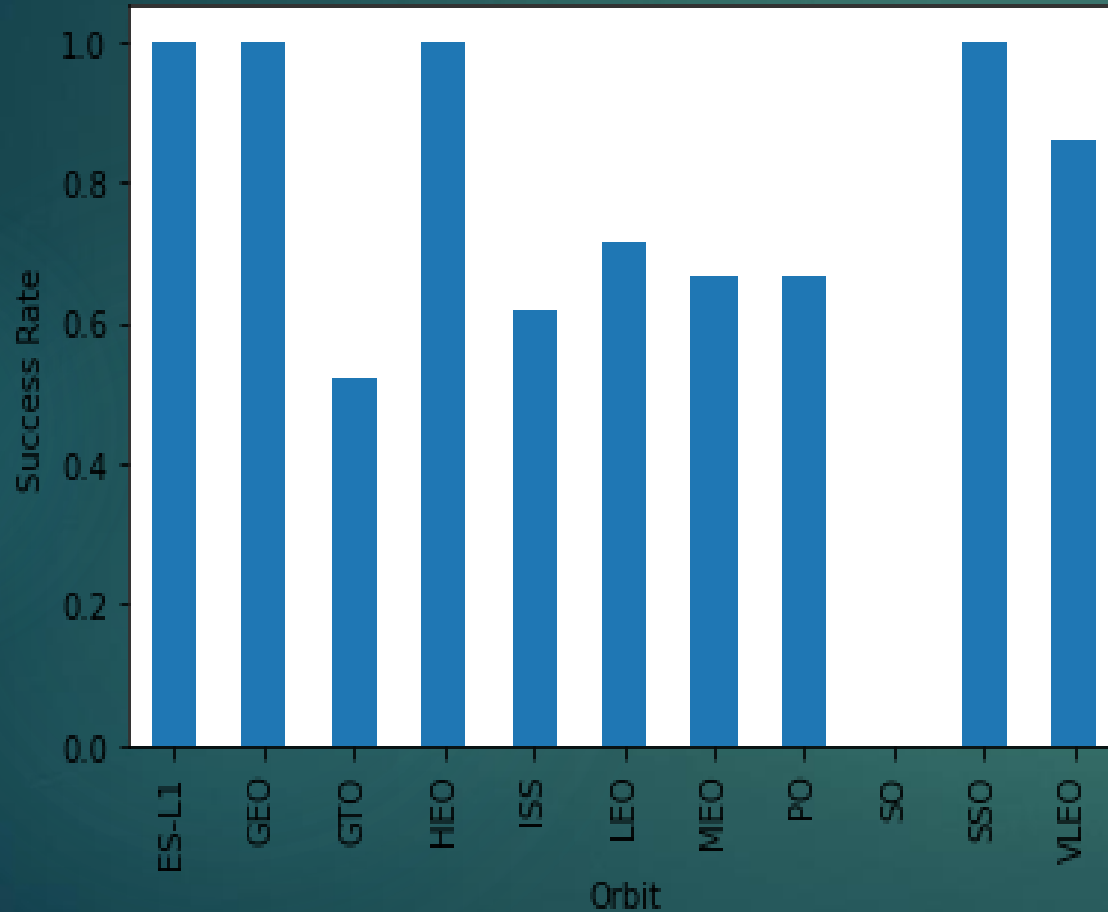
EDA with Visualization



Flight Number Vs Orbit

In this scattered chart of number of flights vs orbit we can see that In the LEO orbit the Success appears to be related to the number of flights; on the other hand, there seem to be no relationship between flight number when in GTO orbit.

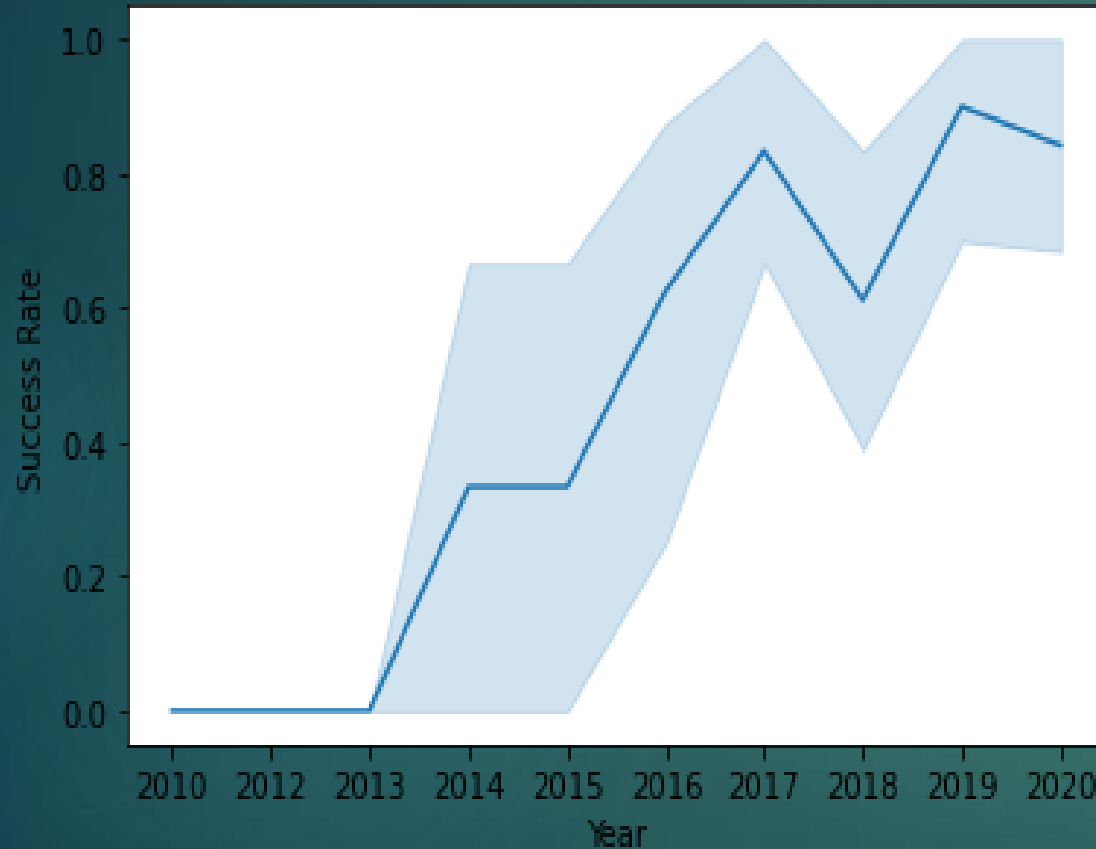
EDA with Visualization



Success Rate Vs Orbit

This bar chart shows that launches aimed at ES-L1, GEO, HEO, and SSO, have a higher success rate than most with SO orbit having the lowest success rate

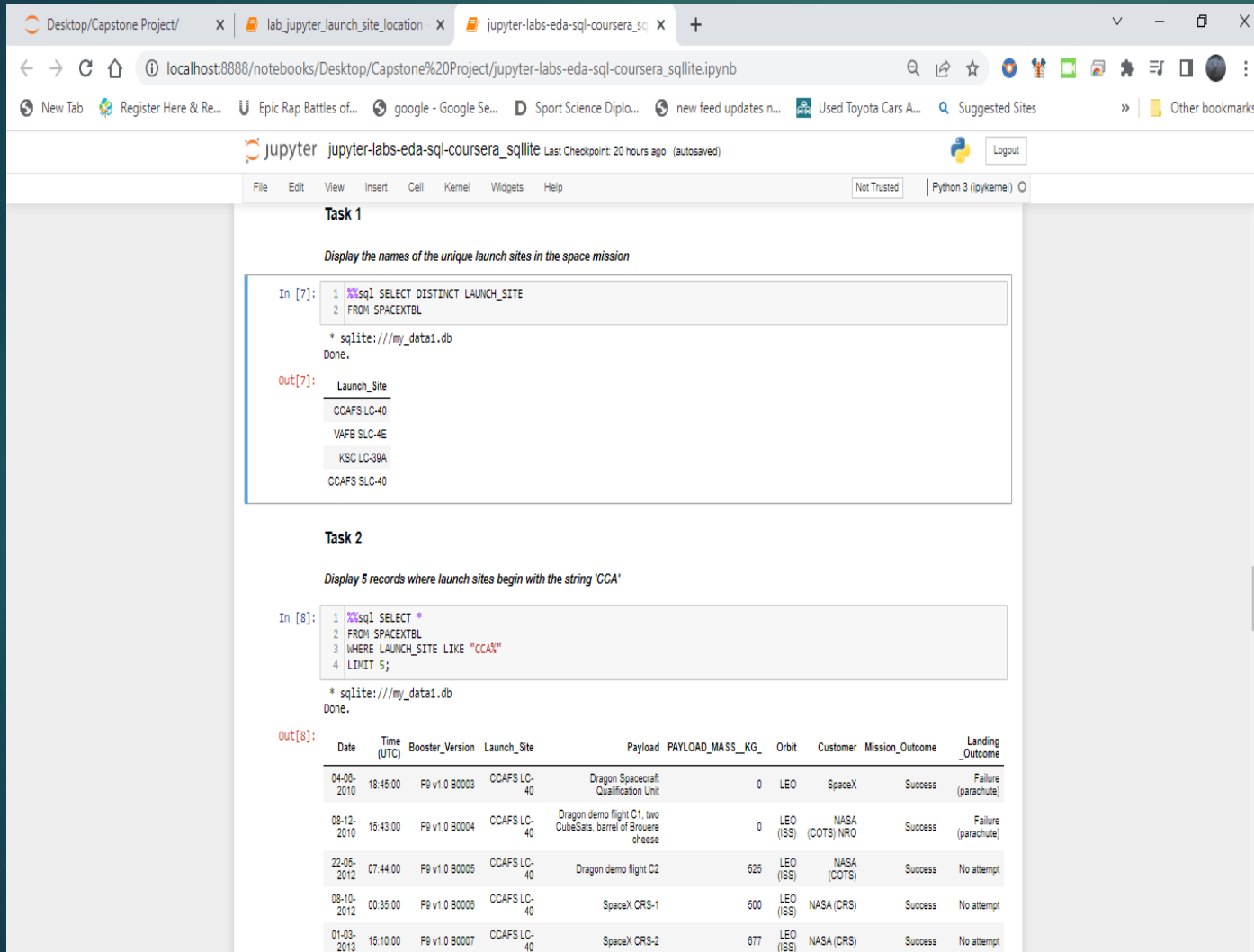
EDA with Visualization



Success Rate Vs Year

In this chart we can observe that the success rate since 2013 kept increasing until a dip in 2018, and a spike from there till 2020.

EDA with SQL



Task 1

Display the names of the unique launch sites in the space mission

```
In [7]: 1 %sql SELECT DISTINCT LAUNCH_SITE
        2 FROM SPACEXTBL

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

Out[7]:

Launch_Site
CCAFS LC-40
VAFB SLC-4E
KSC LC-39A
CCAFS SLC-40

Task 2

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

```
In [8]: 1 %sql SELECT *
        2 FROM SPACEXTBL
        3 WHERE LAUNCH_SITE LIKE "CCA%"
        4 LIMIT 5;

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

Out[8]:

Date	Time (UTC)	Booster_Version	Launch_Site	Payload	PAYLOAD_MASS_KG	Orbit	Customer	Mission_Outcome	Landing_Outcome
04-06-2010	18:45:00	F9 v1.0 B0003	CCAFS LC-40	Dragon Spacecraft Qualification Unit	0	LEO	SpaceX	Success	Failure (parachute)
08-12-2010	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)
22-05-2012	07:44:00	F9 v1.0 B0005	CCAFS LC-40	Dragon demo flight C2	525	LEO (ISS)	NASA (COTS)	Success	No attempt
08-10-2012	00:35:00	F9 v1.0 B0006	CCAFS LC-40	SpaceX CRS-1	500	LEO (ISS)	NASA (CRS)	Success	No attempt
01-03-2013	15:10:00	F9 v1.0 B0007	CCAFS LC-40	SpaceX CRS-2	677	LEO (ISS)	NASA (CRS)	Success	No attempt

- Display the unique launch sites names
- 5 records where Launch sites names begins with "CCA"

EDA with SQL

Desktop/Capstone Project/ x lab_jupyter_launch_site_location x jupyter-labs-eda-sql-coursera_sq x +

localhost:8888/notebooks/Desktop/Capstone%20Project/jupyter-labs-eda-sql-coursera_sqlite.ipynb

New Tab Register Here & Re... Epic Rap Battles of... google - Google Se... Sport Science Diplo... new feed updates n... Used Toyota Cars A... Suggested Sites Other bookmarks

jupyter jupyter-labs-eda-sql-coursera_sqlite Last Checkpoint: 20 hours ago (autosaved) Logout

File Edit View Insert Cell Kernel Widgets Help Not Trusted Python 3 (pykernel)

Task 3

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

In [9]:

```
1 %sql SELECT SUM(PAYLOAD_MASS_KG_)
2 FROM SPACEXTBL
3 WHERE CUSTOMER = "NASA (CRS)";
```

* sqlite:///my_data1.db
Done.

Out[9]:

SUM(PAYLOAD_MASS_KG_)
45596

Task 4

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

In [10]:

```
1 %sql
2 SELECT AVG(PAYLOAD_MASS_KG_)
3 FROM SPACEXTBL
4 WHERE BOOSTER_VERSION = "F9 v1.1";
```

* sqlite:///my_data1.db
Done.

Out[10]:

AVG(PAYLOAD_MASS_KG_)
2928.4

Task 5

List the date when the first successful landing outcome in ground pad was achieved.
Hint: Use min function

In [38]:

```
1 %sql
2 SELECT MIN("Date"), "Landing_Outcomes"
3 FROM SPACEXTBL
4 WHERE "Landing_Outcomes" = "Success (ground pad)";
```

* sqlite:///my_data1.db
Done.

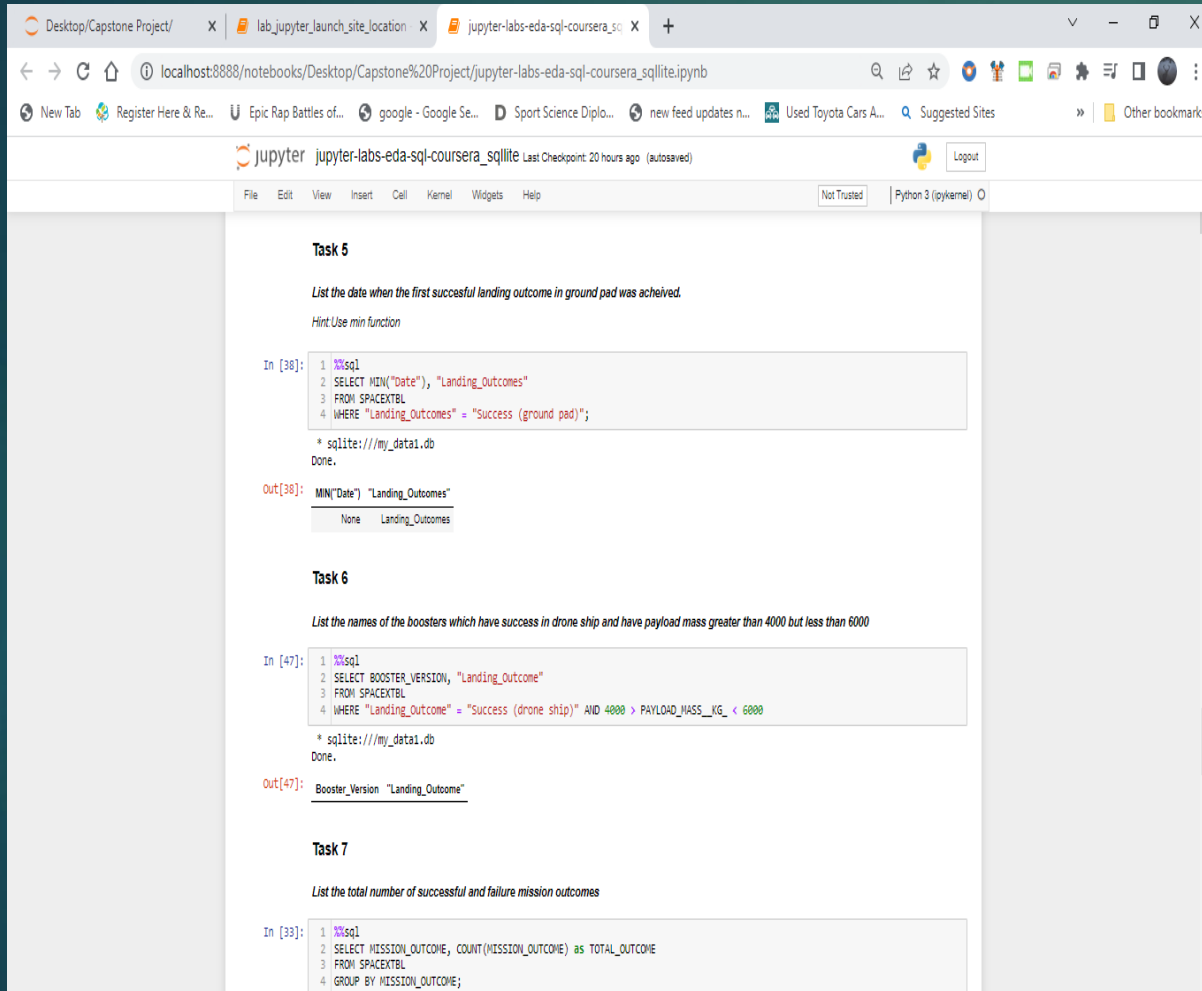
Out[38]:

MIN("Date")	"Landing_Outcomes"
None	Landing_Outcomes

- Total Payload carried by boosters launched by NASA (CRS)
- Average payload carried by booster version F9 v1.1
- Date when the first successful landing outcome in ground pad was achieved

EDA with SQL

- Names of booster with success in drone ship having payload mass greater than 4000 and less than 6000



The screenshot shows a JupyterLab environment with a browser window displaying a Jupyter notebook. The notebook contains three tasks, each with a description, a hint, and a SQL query. The first task, 'Task 5', asks for the date when the first successful landing outcome in ground pad was achieved. The second task, 'Task 6', asks for the names of boosters which have success in drone ship and have payload mass greater than 4000 but less than 6000. The third task, 'Task 7', asks for the total number of successful and failure mission outcomes. The SQL queries are written in a code cell and the results are displayed in a table format.

```
Task 5

List the date when the first successful landing outcome in ground pad was achieved.

Hint Use min function

In [38]: %sql
1 SELECT MIN("Date"), "Landing_Outcomes"
2 FROM SPACEXTBL
3 WHERE "Landing_Outcomes" = "Success (ground pad)";
* sqlite:///my_data1.db
Done.

Out[38]: MIN("Date") "Landing_Outcomes"
None Landing_Outcomes

Task 6

List the names of the boosters which have success in drone ship and have payload mass greater than 4000 but less than 6000

In [47]: %sql
1 SELECT BOOSTER_VERSION, "Landing_Outcome"
2 FROM SPACEXTBL
3 WHERE "Landing_Outcome" = "Success (drone ship)" AND 4000 < PAYLOAD_MASS_KG < 6000
* sqlite:///my_data1.db
Done.

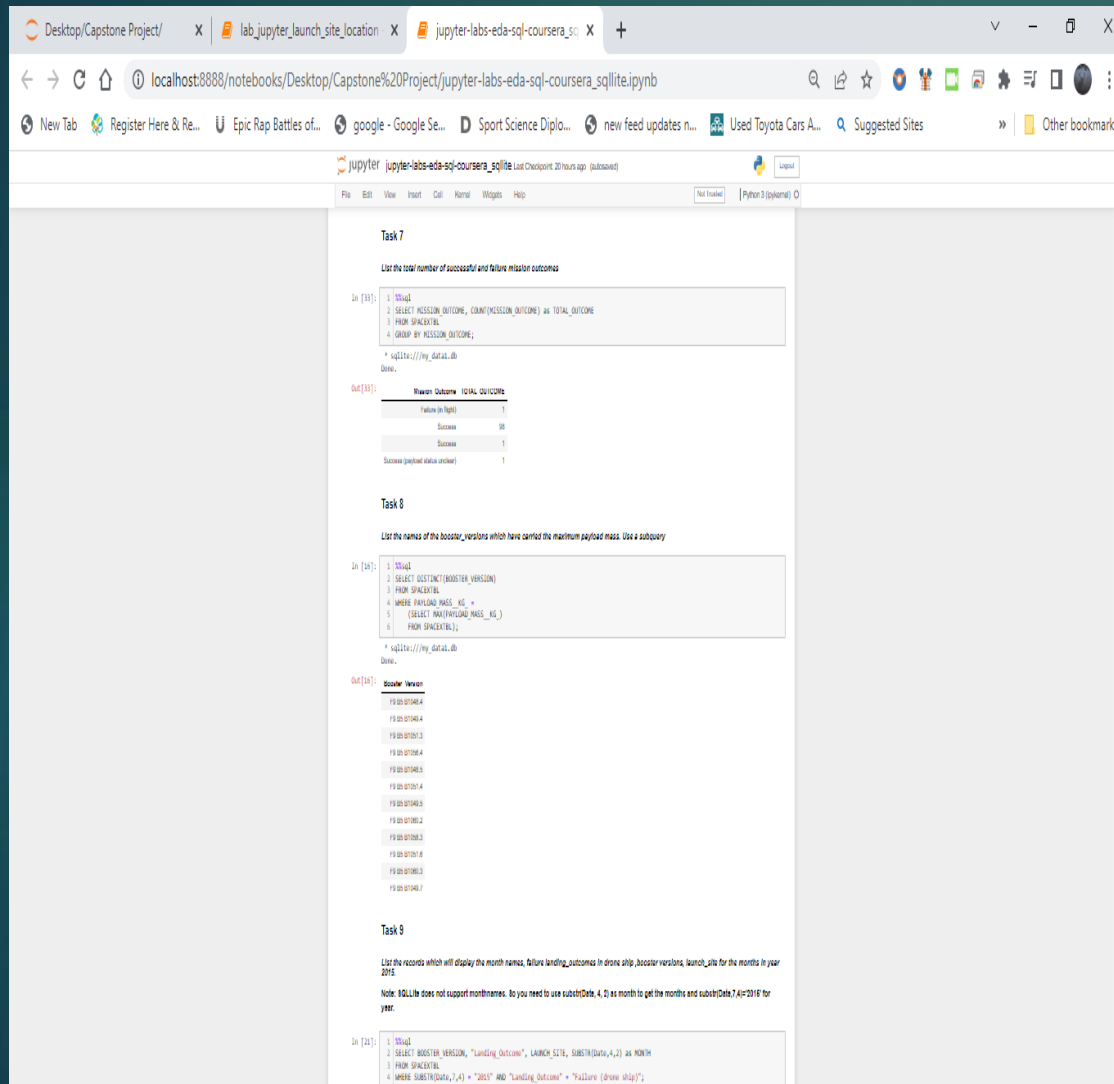
Out[47]: Booster_Version "Landing_Outcome"

Task 7

List the total number of successful and failure mission outcomes

In [33]: %sql
1 SELECT MISSION_OUTCOME, COUNT(MISSION_OUTCOME) as TOTAL_OUTCOME
2 FROM SPACEXTBL
3 GROUP BY MISSION_OUTCOME;
```


EDA with SQL



The screenshot shows a Jupyter Notebook with three tasks. Task 7 asks for the total number of successful and failure mission outcomes. The SQL query uses a subquery to count outcomes and group by mission outcome. The result is a table with two rows: Failure (3) and Success (10). Task 8 asks for the names of booster versions which have carried the maximum payload mass. The SQL query uses a subquery to find the maximum payload mass and then selects booster versions that match. The result is a list of 15 booster versions. Task 9 asks for the records which will display the month names, failure landing outcomes in drone ship, booster versions, launch site for the months in year 2015. The SQL query uses a subquery to filter for the year 2015 and then selects the required columns. The result is a list of 15 records.

Task 7

List the total number of successful and failure mission outcomes

```
In [18]: 1. SELECT
2. SELECT MISSION_OUTCOME, COUNT(MISSION_OUTCOME) AS TOTAL_OUTCOME
3. FROM SPACEXTABLE
4. GROUP BY MISSION_OUTCOME;
```

Out[18]:

Mission Outcome	TOTAL_OUTCOME
Failure (in flight)	3
Success	10
Success (partial static checkout)	1

Task 8

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

```
In [19]: 1. SELECT
2. SELECT DISTINCT(BOOSTER_VERSION)
3. FROM SPACEXTABLE
4. WHERE PAYLOAD_MASS_KG =
5. (SELECT MAX(PAYLOAD_MASS_KG)
6. FROM SPACEXTABLE);
```

Out[19]:

Booster Version
FD 00 01040.4
FD 00 01040.4
FD 00 01040.3
FD 00 01040.4
FD 00 01040.3
FD 00 01040.4
FD 00 01040.3
FD 00 01040.3
FD 00 01040.2
FD 00 01040.3
FD 00 01040.4
FD 00 01040.3
FD 00 01040.7

Task 9

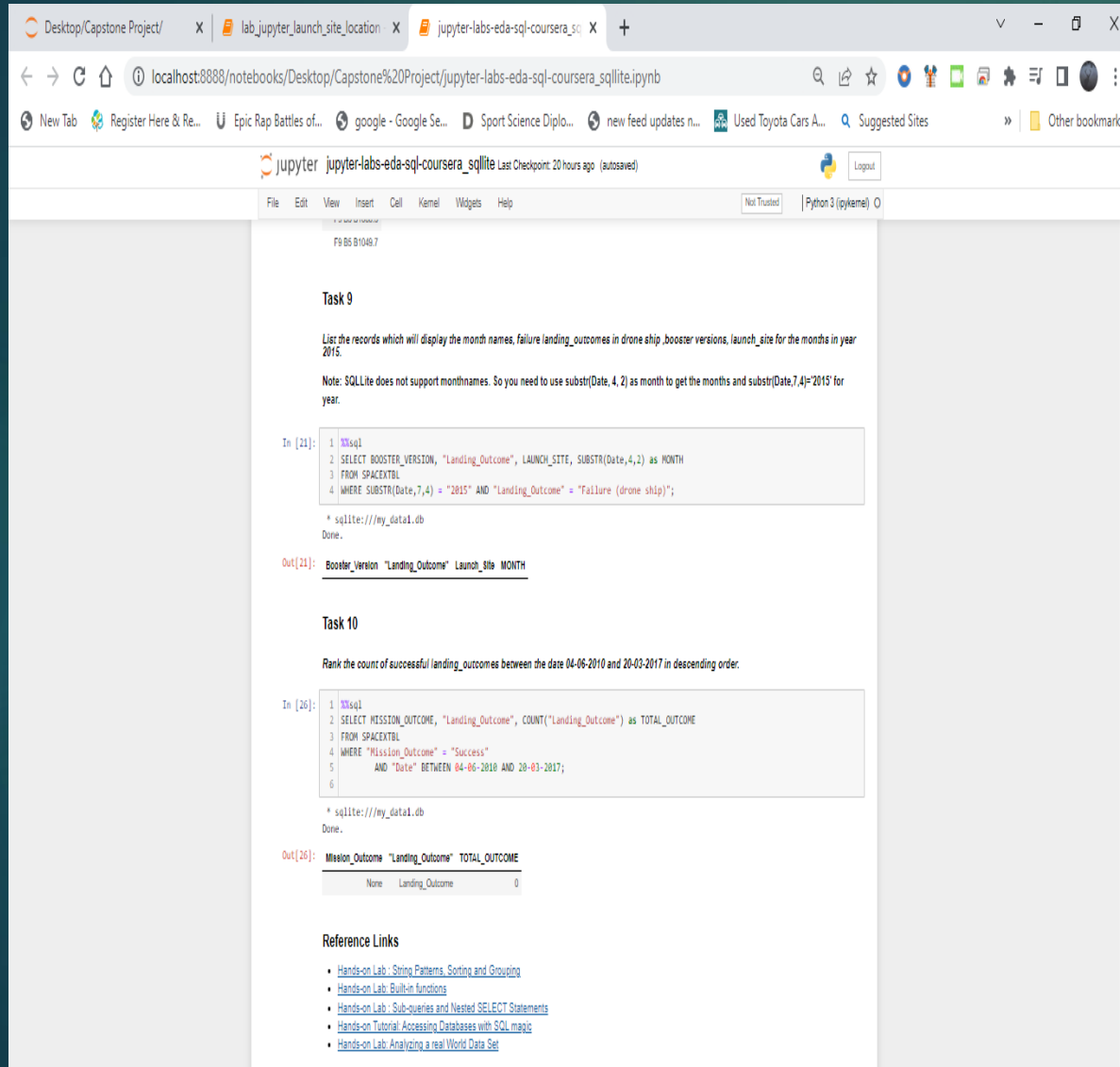
List the records which will display the month names, failure landing_outcomes in drone ship, booster versions, launch_site for the months in year 2015

Note: SQLLite does not support monthnames. So you need to use substr(Date, 4, 2) as month to get the months and substr(Date, 7, 4) as year.

```
In [20]: 1. SELECT
2. SELECT BOOSTER_VERSION, "Landing_Outcome", LAUNCH_SITE, SUBSTR(Date, 4, 2) AS MONTH
3. FROM SPACEXTABLE
4. WHERE SUBSTR(Date, 7, 4) = "2015" AND "Landing_Outcome" = "Failure (drone ship)";
```

- List of total number of successful and failure mission outcome.
- Names of booster versions which have carried the maximum payload mass.

EDA with SQL



The screenshot shows a JupyterLab interface with two tasks. Task 9 involves a SQL query to list records for the year 2015, showing month names, failure landing outcomes, drone ship names, and booster versions. Task 10 involves a SQL query to rank the count of successful landing outcomes between 04-06-2010 and 20-03-2017 in descending order.

Task 9

List the records which will display the month names, failure landing_outcomes in drone ship, booster versions, launch_site for the months in year 2015.

Note: SQLite does not support monthnames. So you need to use `substr(Date, 4, 2)` as month to get the months and `substr(Date, 7, 4)` for year.

```
In [21]: 1 %sql
2 SELECT BOOSTER_VERSION, "Landing_Outcome", LAUNCH_SITE, SUBSTR(Date, 4, 2) as MONTH
3 FROM SPACEXTBL
4 WHERE SUBSTR(Date, 7, 4) = "2015" AND "Landing_Outcome" = "Failure (drone ship)";
```

* sqlite:///my_data1.db
Done.

Out[21]:

Booster_Version	"Landing_Outcome"	Launch_Site	MONTH
-----------------	-------------------	-------------	-------

Task 10

Rank the count of successful landing_outcomes between the date 04-06-2010 and 20-03-2017 in descending order.

```
In [26]: 1 %sql
2 SELECT MISSION_OUTCOME, "Landing_Outcome", COUNT("Landing_Outcome") as TOTAL_OUTCOME
3 FROM SPACEXTBL
4 WHERE "Mission_Outcome" = "Success"
5 AND "Date" BETWEEN 04-06-2010 AND 20-03-2017;
```

* sqlite:///my_data1.db
Done.

Out[26]:

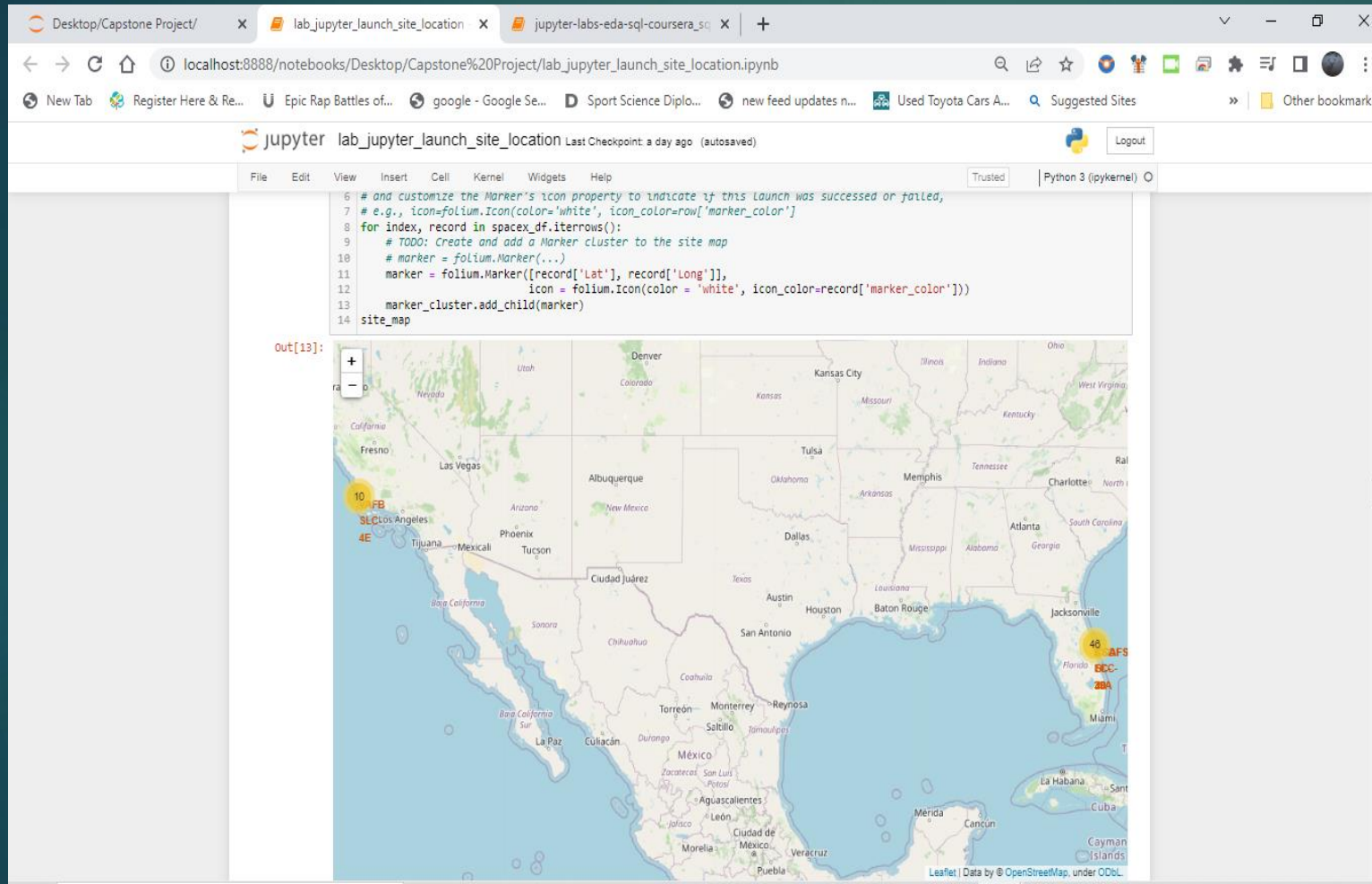
Mission_Outcome	"Landing_Outcome"	TOTAL_OUTCOME
None	Landing_Outcome	0

Reference Links

- [Hands-on Lab: String Patterns, Sorting and Grouping](#)
- [Hands-on Lab: Built-in functions](#)
- [Hands-on Lab: Sub-queries and Nested SELECT Statements](#)
- [Hands-on Tutorial: Accessing Databases with SQL magic](#)
- [Hands-on Lab: Analyzing a real World Data Set](#)

- Records which will display the names, failure landing outcome in drone ship, booster versions launch_size for the months in year 2015
- Rank the count of successful landing outcomes between the date 04-06-2010 and 20-03-2017 in descending order

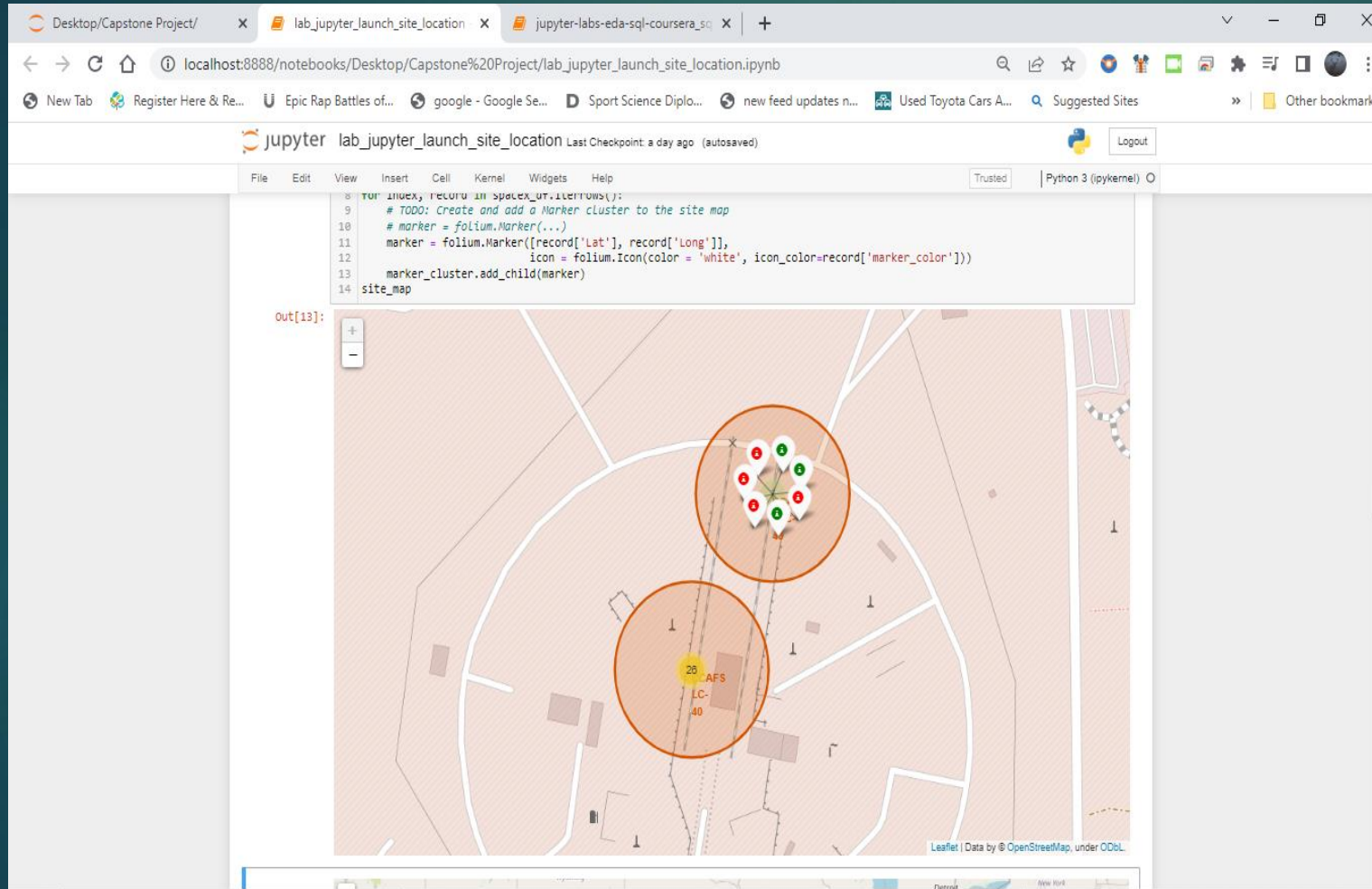
Interactive Map with folium



Folium circle and label for each launch site on the map.

This shows the major coordinates on the map where launches took place.

Interactive Map with folium



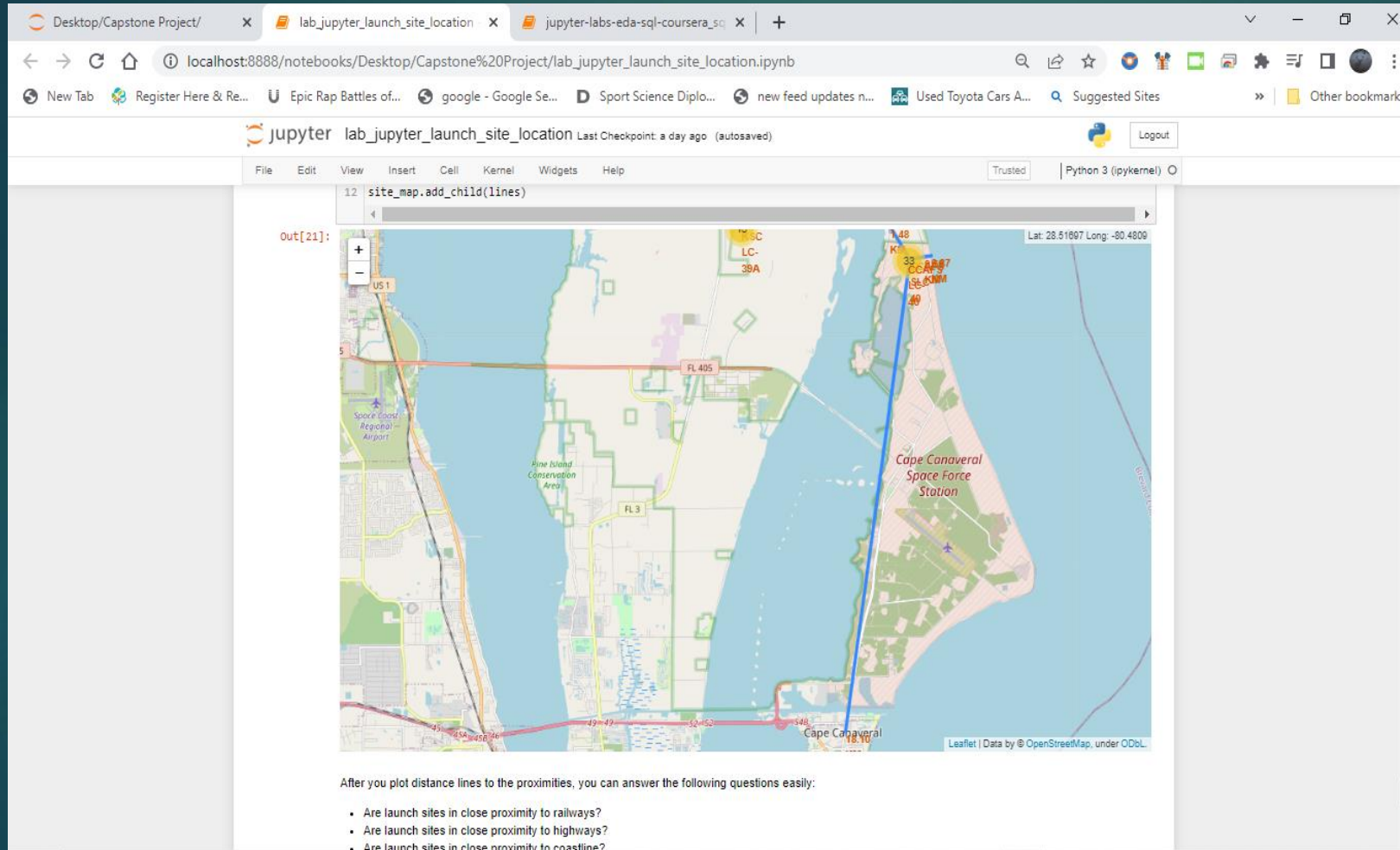
Marker cluster for each launch result

This shows the various launches successful (green), failed (red) at CCAFS SLC-40 launch site

Interactive Map with folium

Proximity Map

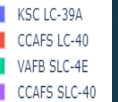
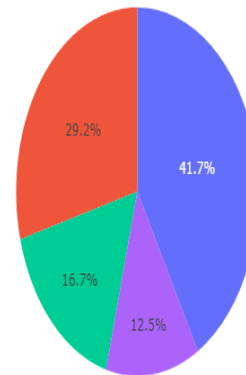
This map shows lines with distance from launch site to major city, highway, coastline, and railway.



Interactive Visual with Plotly Dash

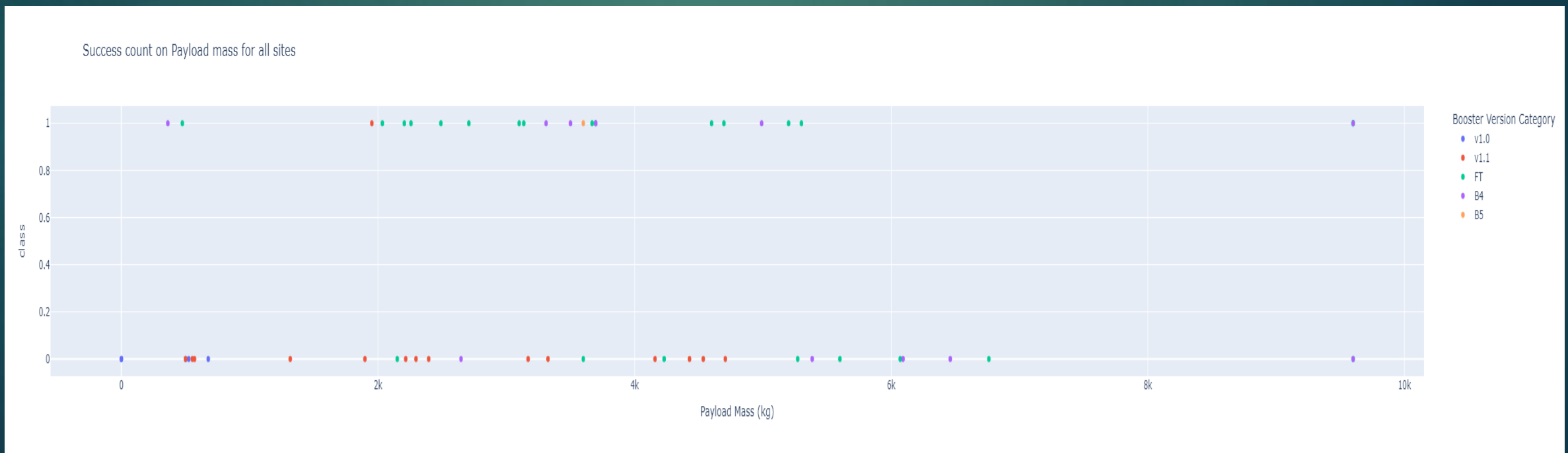
The pie chart shows that KSC LC-39A has the largest successful launches with a 41% launch success count

Success Count for all launch sites



Interactive Visual with Plotly Dash

- The scattered plot shows that payload within 2k and 4k range has the highest success range
- Payload within 6k and 8k range has the lowest success rate
- F9 Booster Version FT has the highest launch success rate





Predictive analysis Result Accuracy and confusion matrix

Splitting our data into train and test data set.

[illegible]

```

TASK 3

Use the function train_test_split to split the data X and Y into training and test data. Set the parameter test_size to 0.2 and random_state to 2. The training data and test data should be assigned to the following labels

X_train, X_test, Y_train, Y_test

In [8]:
1 X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.2, random_state=2)

we can see we only have 18 test samples.

In [9]:
1 Y_test.shape

Out[9]: (18,)
```

Logistic Regression

Creation of a LR object, fit the object to find the best parameter from the dictionary parameters.

Calculating the accuracy, and plotting the confusion matrix

TASK 4

Create a logistic regression object then create a GridSearchCV object `logreg_cv` with `cv = 10`. Fit the object to find the best parameters from the dictionary parameters.

In [10]:

```
1 parameters = {'C': [0.01, 0.1, 1],  
2               'penalty': ['l2'],  
3               'solver': ['lbfgs']}
```

In [11]:

```
1 parameters = {'C': [0.01, 0.1, 1], 'penalty': ['l2'], 'solver': ['lbfgs']} # l1 Lasso l2 ridge  
2 lr = LogisticRegression()  
3 logreg_cv = GridSearchCV(lr, parameters, cv=10)  
4 logreg_cv.fit(X_train, Y_train)  
5
```

cision. If you wish to review your current use, check the release note link for additional information.
Deprecated in NumPy 1.20; for more details and guidance: <https://numpy.org/devdocs/release/1.20.0-notes.html#deprecations>
dtype=np.int)
/home/jupyterlab/conda/envs/python/lib/python3.7/site-packages/sklearn/model_selection/_search.py:841: DeprecationWarning:
The default of the 'iid' parameter will change from True to False in version 0.22 and will be removed in 0.24. This will change
numeric results when test-set sizes are unequal.
DeprecationWarning:
/home/jupyterlab/conda/envs/python/lib/python3.7/site-packages/sklearn/utils/fixes.py:357: DeprecationWarning: distutils version
classes are deprecated. Use packaging.version instead.
if _joblib__version__ > LooseVersion("0.12"):

```
Out[11]: GridSearchCV(cv=10, error_score='raise-deprecated',  
estimator=LogisticRegression(C=1.0, class_weight=None, dual=False, fit_intercept=True,  
intercept_scaling=1, max_iter=100, multi_class='warn',  
n_jobs=None, penalty='l2', random_state=None, solver='warn',  
tol=0.0001, verbose=0, warn_start=False),  
fit_params=None, iid='warn', n_jobs=None,  
param_grid={'C': [0.01, 0.1, 1], 'penalty': ['l2'], 'solver': ['lbfgs']},  
pre_dispatch='2*n_jobs', refit=True, return_train_score='warn',  
scoring=None, verbose=0)
```

We output the `GridSearchCV` object for logistic regression. We display the best parameters using the data attribute `best_params_` and the accuracy on the validation data using the data attribute `best_score_`.

In [12]:

```
1 print("tuned hyperparameters : (best parameters) ", logreg_cv.best_params_)  
2 print("accuracy : ", logreg_cv.best_score_)
```

tuned hyperparameters : (best parameters) {'C': 0.01, 'penalty': 'l2', 'solver': 'lbfgs'}
accuracy : 0.8472222222222222

TASK 5

Calculate the accuracy on the test data using the method `score`:

In [13]:

```
1 print("Test Accuracy:", logreg_cv.score(X_test, Y_test))
```

Test Accuracy: 0.8333333333333334

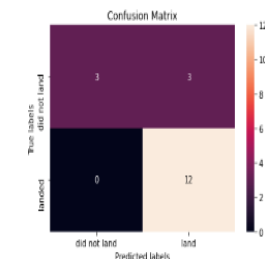
/home/jupyterlab/conda/envs/python/lib/python3.7/site-packages/sklearn/linear_model/base.py:283: DeprecationWarning: 'np.int' is a deprecated alias for the builtin 'int'. To silence this warning, use 'int' by itself. Doing this will not modify any behavior and is safe. When replacing 'np.int', you may wish to use e.g. 'np.int64' or 'np.int32' to specify the precision. If you wish to review your current use, check the release note link for additional information.
Deprecated in NumPy 1.20; for more details and guidance: <https://numpy.org/devdocs/release/1.20.0-notes.html#deprecations>
Indices = (scores > 0).astype(np.int)

Lets look at the confusion matrix:

In [14]:

```
1 yhat = logreg_cv.predict(X_test)  
2 plot_confusion_matrix(Y_test, yhat)
```

/home/jupyterlab/conda/envs/python/lib/python3.7/site-packages/sklearn/linear_model/base.py:283: DeprecationWarning: 'np.int' is a deprecated alias for the builtin 'int'. To silence this warning, use 'int' by itself. Doing this will not modify any behavior and is safe. When replacing 'np.int', you may wish to use e.g. 'np.int64' or 'np.int32' to specify the precision. If you wish to review your current use, check the release note link for additional information.
Deprecated in NumPy 1.20; for more details and guidance: <https://numpy.org/devdocs/release/1.20.0-notes.html#deprecations>
Indices = (scores > 0).astype(np.int)



Support Vector Machine

Create a SVM object, fit and search for the best parameter from the parameter dictionary

Calculate accuracy and plot confusion matrix

TASK 6

Create a support vector machine object then create a `GridSearchCV` object `svm_cv` with `cv=10`. Fit the object to find the best parameters from the dictionary `parameters`.

```
In [15]: parameters = {'kernel': ('linear', 'rbf', 'poly', 'rbf', 'sigmoid'),
1         'C': np.logspace(-3, 3, 5),
2         'gamma': np.logspace(-3, 3, 5)}
3 svm = SVC()
```

```
In [16]: 1 svm_cv=GridSearchCV(svm, parameters,cv=10)
2 svm_cv.fit(X_train,Y_train)
```

```
DeprecationWarning: np.int is a deprecated alias for the builtin 'int'. To silence this warning, use 'int' by itself. Doing this will not modify any behavior and is safe. When replacing 'np.int', you may wish to use e.g. 'np.int64' or 'np.int32' to specify the precision. If you wish to review your current use, check the release note link for additional information.
DeprecationWarning: The default of the 'iid' parameter will change from True to False in version 0.22 and will be removed in 0.24. This will change numeric results when test-set sizes are unequal.
DeprecationWarning: dtype=np.int)
```

```
Out[16]: GridSearchCV(cv=10, error_score='raise-deprecating',
    estimator=SVC(C=1.0, cache_size=200, class_weight=None, coef0=0.0,
    decision_function_shape='ovr', degree=3, gamma='auto-deprecated',
    kernel='rbf', max_iter=1, probability=False, random_state=None,
    shrinking=True, tol=0.001, verbose=False),
    fit_params=None, iid='warn', n_jobs=None,
    param_grid={'kernel': ('linear', 'rbf', 'poly', 'rbf', 'sigmoid'), 'C': array([1.00000e-03, 3.16228e-02, 1.00000e+00, 3.16228e+01, 1.00000e+03]), 'gamma': array([1.00000e-03, 3.16228e-02, 1.00000e+00, 3.16228e+01, 1.00000e+03])},
    pre_dispatch='2*n_jobs', refit=True, return_train_score='warn',
    scoring=None, verbose=0)
```

```
In [17]: 1 print("tuned hyperparameters : (best parameters) ",svm_cv.best_params_)
2 print("accuracy : ",svm_cv.best_score_)

tuned hyperparameters : (best parameters) {'C': 1.0, 'gamma': 0.03162277660168379, 'kernel': 'sigmoid'}
accuracy : 0.8472222222222222
```

TASK 7

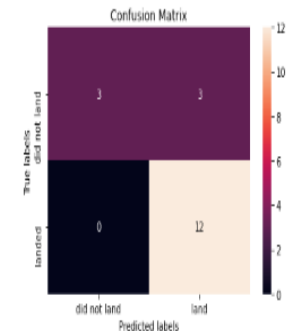
Calculate the accuracy on the test data using the method `score`:

```
In [18]: 1 print("Test Accuracy:", svm_cv.score(X_test,Y_test))
```

Test Accuracy: 0.8333333333333334

We can plot the confusion matrix

```
In [19]: 1 yhat=svm_cv.predict(X_test)
2 plot_confusion_matrix(Y_test,yhat)
```



Decision Tree Classifier

Create a decision tree classifier object, fit to search for the best parameter from the parameter dictionary.

TASK 8

Edit Attachments

Create a decision tree classifier object then create a `GridSearchCV` object `tree_cv` with `cv=10`. Fit the object to find the best parameters from the dictionary parameters.

In [20]:

Edit Attachments

```
1 parameters = {'criterion': ['gini', 'entropy'],
2               'splitter': ['best', 'random'],
3               'max_depth': [2*n for n in range(1,10)],
4               'max_features': ['auto', 'sqrt'],
5               'min_samples_leaf': [1, 2, 4],
6               'min_samples_split': [2, 5, 10]}
7
8 tree = DecisionTreeClassifier()
```

In [21]:

Edit Attachments

```
1 tree_cv = GridSearchCV(tree, parameters, cv=10)
2 tree_cv.fit(X_train, Y_train)
```

DeprecationWarning

/home/jupyterlab/conda/envs/python/lib/python3.7/site-packages/sklearn/tree/tree.py:148: DeprecationWarning: 'np.int' is a deprecated alias for the builtin 'int'. To silence this warning, use 'int' by itself. Doing this will not modify any behavior and is safe. When replacing 'np.int', you may wish to use e.g. 'np.int64' or 'np.int32' to specify the precision. If you wish to review your current use, check the release note link for additional information.
Deprecated in NumPy 1.20; for more details and guidance: https://numpy.org/devdocs/release/1.20.0-notes.html#deprecations
y_encoded = np.zeros(y.shape, dtype=np.int)

Out[21]:

GridSearchCV(cv=10, error_score='raise-deprecating',
estimator=DecisionTreeClassifier(class_weight=None, criterion='gini', max_depth=None,
max_features=None, max_leaf_nodes=None,
min_impurity_decrease=0.0, min_impurity_split=None,
min_samples_leaf=1, min_samples_split=2,
min_weight_fraction_leaf=0.0, presort=False, random_state=None,
splitter='best'),
fit_params=None, iid='warn', n_jobs=None,
param_grid={'criterion': ['gini', 'entropy'], 'splitter': ['best', 'random'], 'max_depth': [2, 4, 6, 8, 10, 12, 14, 1
6, 18], 'max_features': ['auto', 'sqrt'], 'min_samples_leaf': [1, 2, 4], 'min_samples_split': [2, 5, 10]},
pre_dispatch='2*n_jobs', refit=True, return_train_score='warn',
scoring=None, verbose=0)

In [22]:

Edit Attachments

```
1 print("tuned hyperparameters : (best parameters) ", tree_cv.best_params_)
2 print("accuracy : ", tree_cv.best_score_)
```

tuned hyperparameters : (best parameters) {'criterion': 'gini', 'max_depth': 6, 'max_features': 'sqrt', 'min_samples_leaf': 2,
'min_samples_split': 2, 'splitter': 'best'}
accuracy : 0.8888888888888888

Calculate accuracy and plot confusion matrix

TASK 9

Edit Attachments

Calculate the accuracy of `tree_cv` on the test data using the method `score`:

In [23]:

Edit Attachments

```
1 print("Test Accuracy:", tree_cv.score(X_test, Y_test))
```

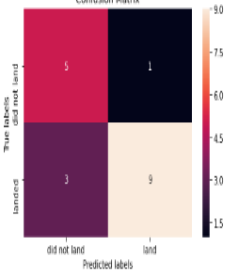
Test Accuracy: 0.7777777777777778

In [24]:

Edit Attachments

```
1 yhat = tree_cv.predict(X_test)
2 plot_confusion_matrix(Y_test, yhat)
```

Confusion Matrix



K Nearest Neighbour

Create a KNN object, fit to find the best parameter from the parameter dictionary.

TASK 10

Edit Attachments

Create a k nearest neighbors object then create a `GridSearchCV` object `knn_cv` with `cv = 10`. Fit the object to find the best parameters from the dictionary `parameters`.

In [25]:

Edit Attachments

```
1 parameters = {'n_neighbors': [1, 2, 3, 4, 5, 6, 7, 8, 9, 10],
2               'algorithm': ['auto', 'ball_tree', 'kd_tree', 'brute'],
3               'p': [1, 2]}
4 knn = KNeighborsClassifier()
```

In [26]:

Edit Attachments

```
1 knn_cv = GridSearchCV(knn, parameters, cv=10)
2 knn_cv.fit(X_train, Y_train)
```

/home/jupyterlab/conda/envs/python/lib/python3.7/site-packages/sklearn/model_selection/_search.py:841: DeprecationWarning: The default of the 'iid' parameter will change from True to False in version 0.22 and will be removed in 0.24. This will change numeric results when test-set sizes are unequal.
 DeprecationWarning)
/home/jupyterlab/conda/envs/python/lib/python3.7/site-packages/sklearn/neighbors/base.py:907: DeprecationWarning: 'np.int' is a deprecated alias for the builtin 'int'. To silence this warning, use 'int' by itself. Doing this will not modify any behavior and is safe. When replacing 'np.int', you may wish to use e.g. 'np.int64' or 'np.int32' to specify the precision. If you wish to review your current use, check the release note link for additional information.
Deprecated In NumPy 1.20; for more details and guidance: https://numpy.org/devdocs/release/1.20.0-notes.html#deprecations
self._y = np.empty(y.shape, dtype=np.int)

Out[26]:

GridSearchCV(cv=10, error_score='raise-deprecating',
 estimator=KNeighborsClassifier(algorithm='auto', leaf_size=30, metric='minkowski',
 metric_params=None, n_jobs=None, n_neighbors=5, p=2,
 weights='uniform'),
 fit_params=None, iid='warn', n_jobs=None,
 param_grid={'n_neighbors': [1, 2, 3, 4, 5, 6, 7, 8, 9, 10], 'algorithm': ['auto', 'ball_tree', 'kd_tree', 'brute'],
 'p': [1, 2]},
 pre_dispatch='2*n_jobs', refit=True, return_train_score='warn',
 scoring=None, verbose=0)

In [27]:

Edit Attachments

```
1 print("tuned hyperparameters : (best parameters) ", knn_cv.best_params_)
2 print("accuracy : ", knn_cv.best_score_)
```

tuned hyperparameters : (best parameters) {'algorithm': 'auto', 'n_neighbors': 9, 'p': 1}
accuracy : 0.8472222222222222

Calculate accuracy on the test data, and plot confusion matrix

TASK 11

Edit Attachments

Calculate the accuracy of `tree_cv` on the test data using the method `score`:

In [28]:

Edit Attachments

```
1 print("Test Accuracy:", knn_cv.score(X_test, Y_test))
```

Test Accuracy: 0.8333333333333334

/home/jupyterlab/conda/envs/python/lib/python3.7/site-packages/sklearn/neighbors/base.py:442: DeprecationWarning: distutils Version classes are deprecated. Use packaging.version instead.
 old_joblib = LooseVersion(joblib_version) < LooseVersion('0.12')
/home/jupyterlab/conda/envs/python/lib/python3.7/site-packages/sklearn/neighbors/base.py:442: DeprecationWarning: distutils Version classes are deprecated. Use packaging.version instead.
 old_joblib = LooseVersion(joblib_version) < LooseVersion('0.12')

We can plot the confusion matrix

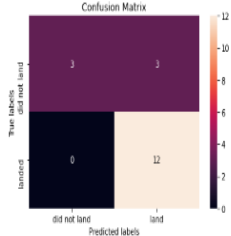
In [29]:

Edit Attachments

```
1 yhat = knn_cv.predict(X_test)
2 plot_confusion_matrix(Y_test, yhat)
```

/home/jupyterlab/conda/envs/python/lib/python3.7/site-packages/sklearn/neighbors/base.py:442: DeprecationWarning: distutils Version classes are deprecated. Use packaging.version instead.
 old_joblib = LooseVersion(joblib_version) < LooseVersion('0.12')
/home/jupyterlab/conda/envs/python/lib/python3.7/site-packages/sklearn/neighbors/base.py:442: DeprecationWarning: distutils Version classes are deprecated. Use packaging.version instead.
 old_joblib = LooseVersion(joblib_version) < LooseVersion('0.12')

Confusion Matrix



CONCLUSION

- Irrespective of the launch site rockets with pay load within 4000kg and 6000kg had a higher failure rate
- Launch success rate climbed as from 2013
- Launch sites has lowest proximity to cities than any Highways, Rail tracks and even coastlines
- KSC LC-39A has the highest success launch rate than any other site
- Orbit SO has the lowest success rate
- Decision tree classifier is the best machine learning model for the task with a higher classification fit accuracy.

A photograph of a Space Shuttle launching from a launchpad, viewed from across a body of water. The shuttle is ascending vertically, leaving a bright white plume of smoke and fire. A massive, billowing cloud of white smoke and steam rises from the launchpad, partially obscuring the sky. To the right of the launchpad, a tall, slender water tower stands. The foreground shows the calm surface of the water, which reflects the bright light from the shuttle's engines. The sky is a deep blue with some lighter clouds. The text "Thank you" is overlaid in a large, white, sans-serif font on the left side of the image.

Thank you