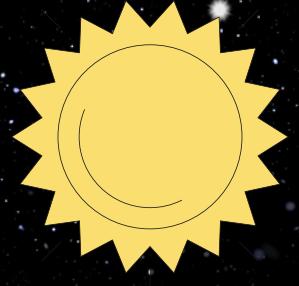


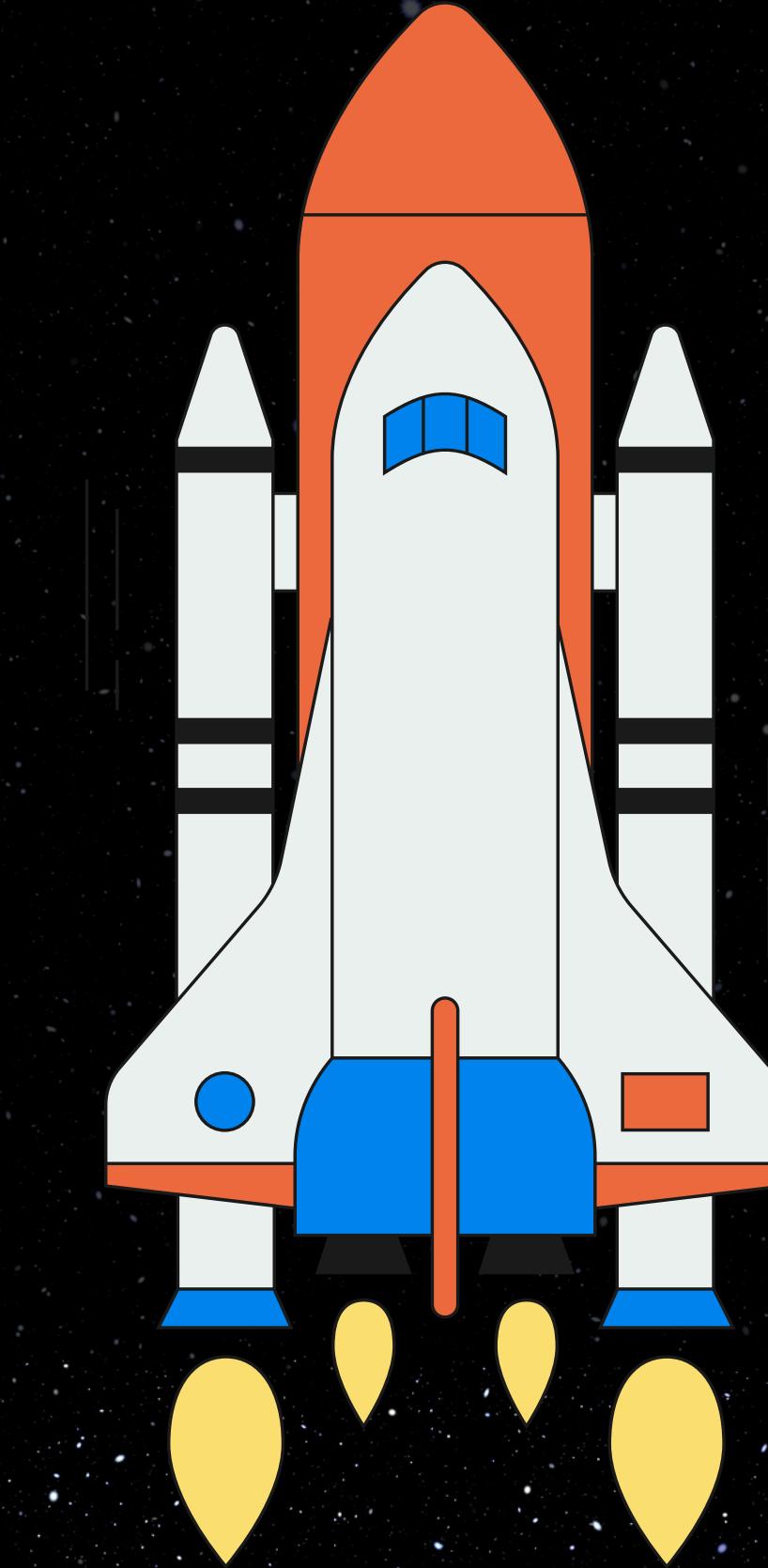
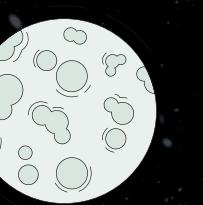
APPLIED DATA SCIENCE CAPSTONE PROJECT

by Nhlanhla Sanele Dludla



Outline

- Executive summary
- Introduction
- Methodology
- Results
- Conclusion

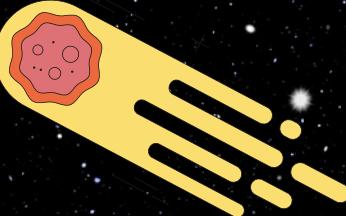


Executive summary

Conventional space exploration is expensive. To cut costs associated with space travel, SpaceX reuses their first stage. In the bid to cut launch costs, in this capstone project we aim to predict whether the first stage of Falcon 9 will land successfully or not on landing pads for reuse. To achieve this, data will be collected cleaned and formatted. In addition exploratory and predictive analyses will be performed. Results will be presented in graph format.

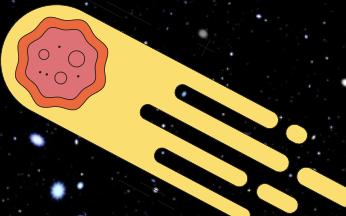
Introduction

Conventional space exploration is expensive. The typical costs to launch space crafts costs up to \$165 Million. To cut costs associated with space travel, SpaceX reuses their first stage. In the bid to cut launch costs, in this capstone project we aim to predict whether the first stage of Falcon 9 will land successfully or not on landing pads for reuse. To achieve this, data will be collected cleaned and formatted. In addition, exploratory and predictive analyses will be performed. Results will be presented in graph format.



Methodology

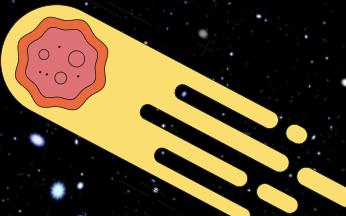
1. Data collection and processing.
 2. Data analyses
 3. Data visualization
 4. Machine learning and prediction
- 



Data collection: API + Data wrangling

1. “<https://api.spacexdata.com/v4/launches/past>” was used to access SpaceX data.
2. Collected data was decoded into a JSON file and IDs were assigned to specific columns.
3. The dataframe was then filtered to include only Falcon 9 launches.
4. `data_falcon9.isnull().sum()` was used to identify missing values.
5. Using `data_falcon9['PayloadMass'] = data_falcon9['PayloadMass'].replace(np.nan, payloadmass_mean)` missing values were replaced with mean values.

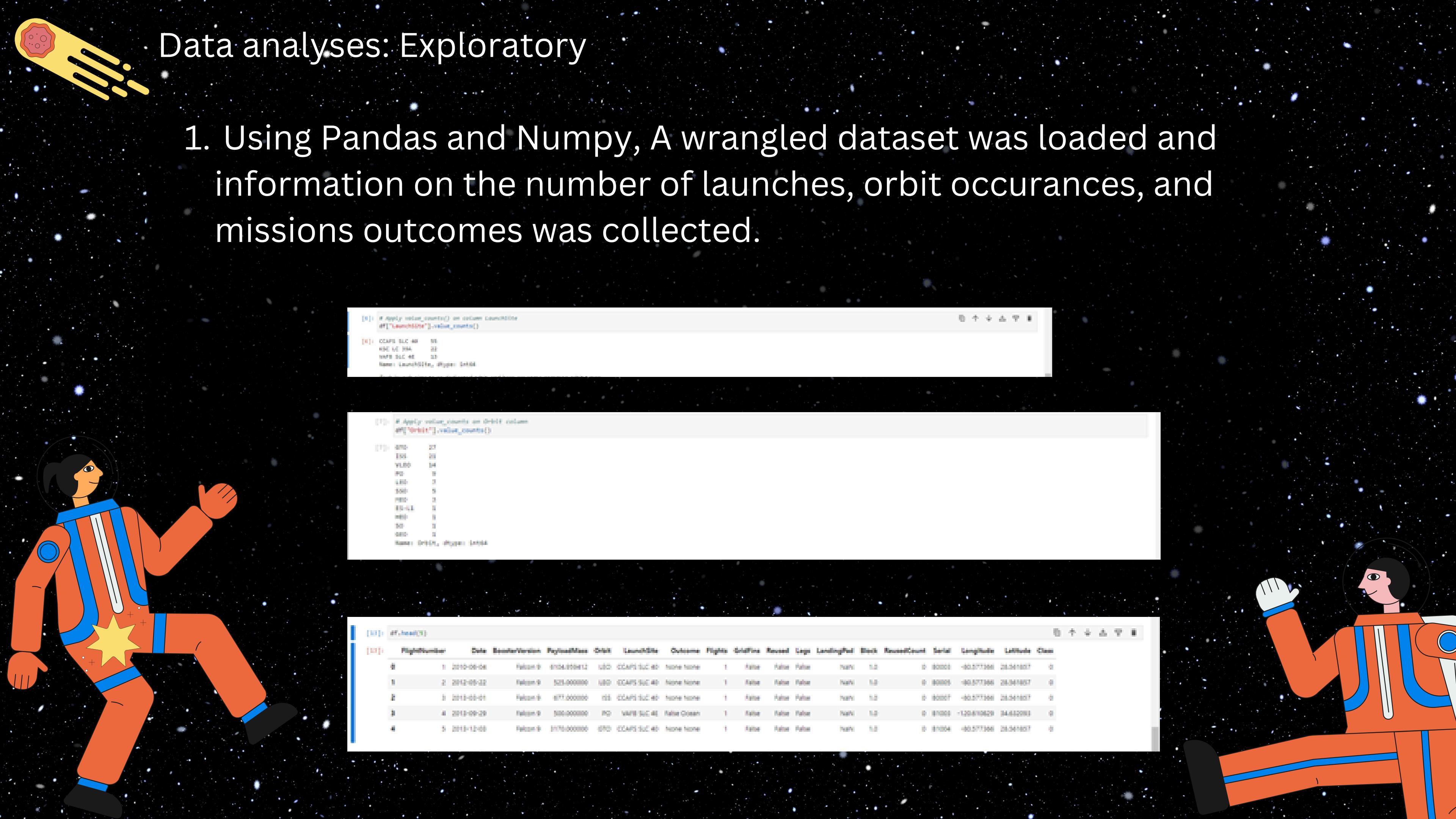
| FlightNumber | Date | BoosterVersion | PayloadMass | Orbit | LaunchSite | Outcome | Flights | GridFins | Reused | Legs | LandingPad | Block | ReusedCount | Serial | Longitude | Latitude |
|--------------|------------|----------------|-------------|-------|--------------|-------------|---------|----------|--------|-------|------------|-------|-------------|--------|-------------|-----------|
| 1 | 2010-06-04 | Falcon 9 | NaN | LEO | CCSFS SLC 40 | None None | 1 | False | False | False | None | 1.0 | 0 | B0003 | -80.577366 | 28.561857 |
| 2 | 2012-05-22 | Falcon 9 | 525.0 | LEO | CCSFS SLC 40 | None None | 1 | False | False | False | None | 1.0 | 0 | B0005 | -80.577366 | 28.561857 |
| 3 | 2013-03-01 | Falcon 9 | 677.0 | ISS | CCSFS SLC 40 | None None | 1 | False | False | False | None | 1.0 | 0 | B0007 | -80.577366 | 28.561857 |
| 4 | 2013-09-29 | Falcon 9 | 500.0 | PO | VAFB SLC 4E | False Ocean | 1 | False | False | False | None | 1.0 | 0 | B1003 | -120.610829 | 34.632093 |
| 5 | 2013-12-03 | Falcon 9 | 3170.0 | GTO | CCSFS SLC 40 | None None | 1 | False | False | False | None | 1.0 | 0 | B1004 | -80.577366 | 28.561857 |

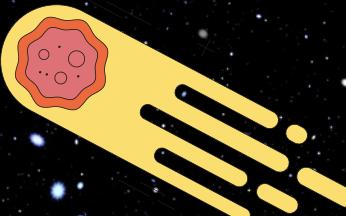


Data collection: Webscrapping

1. Using BeautifulSoup Falcon 9 launch table records were extracted from Wikipedia.
2. HTML tables were parsed and a dataframe was created.

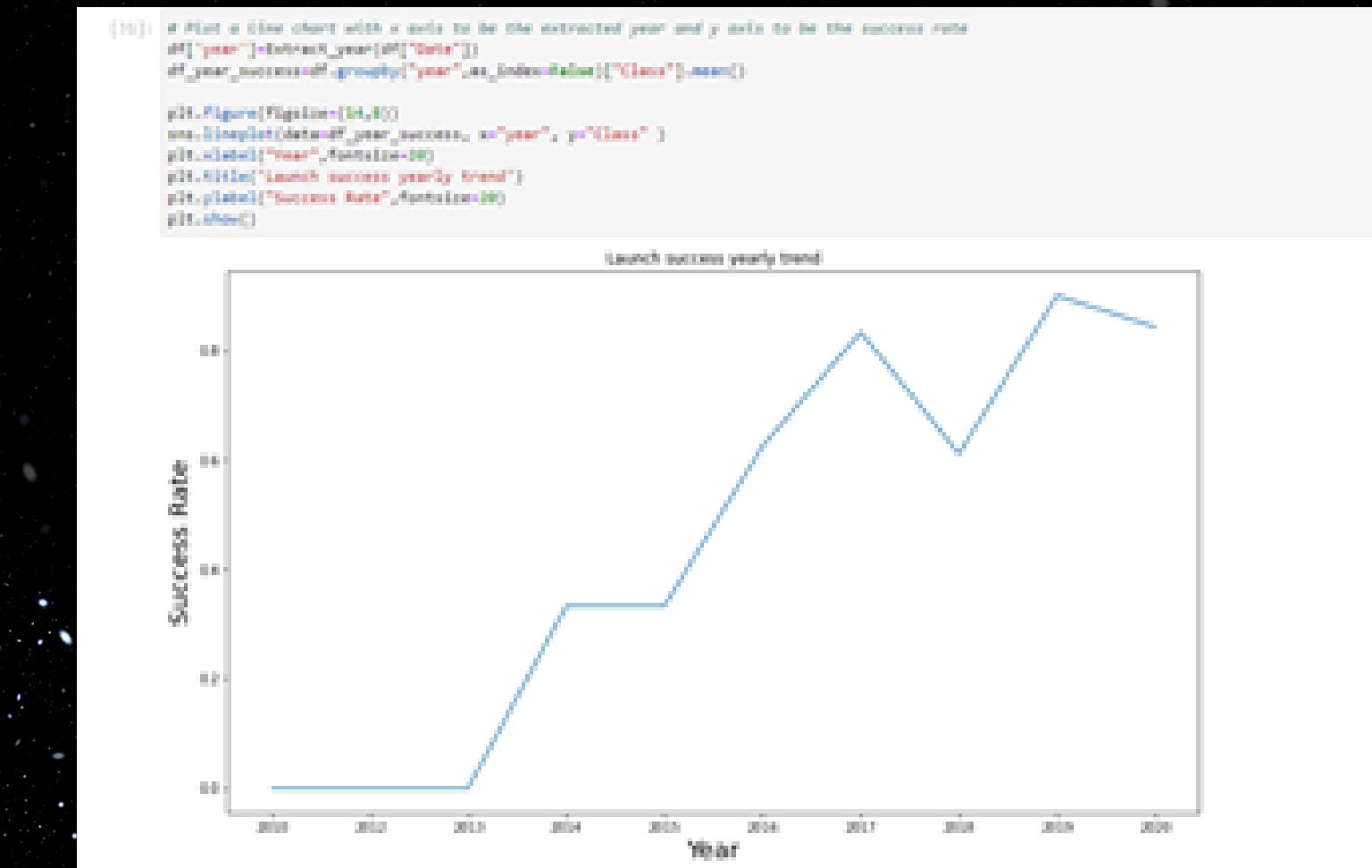
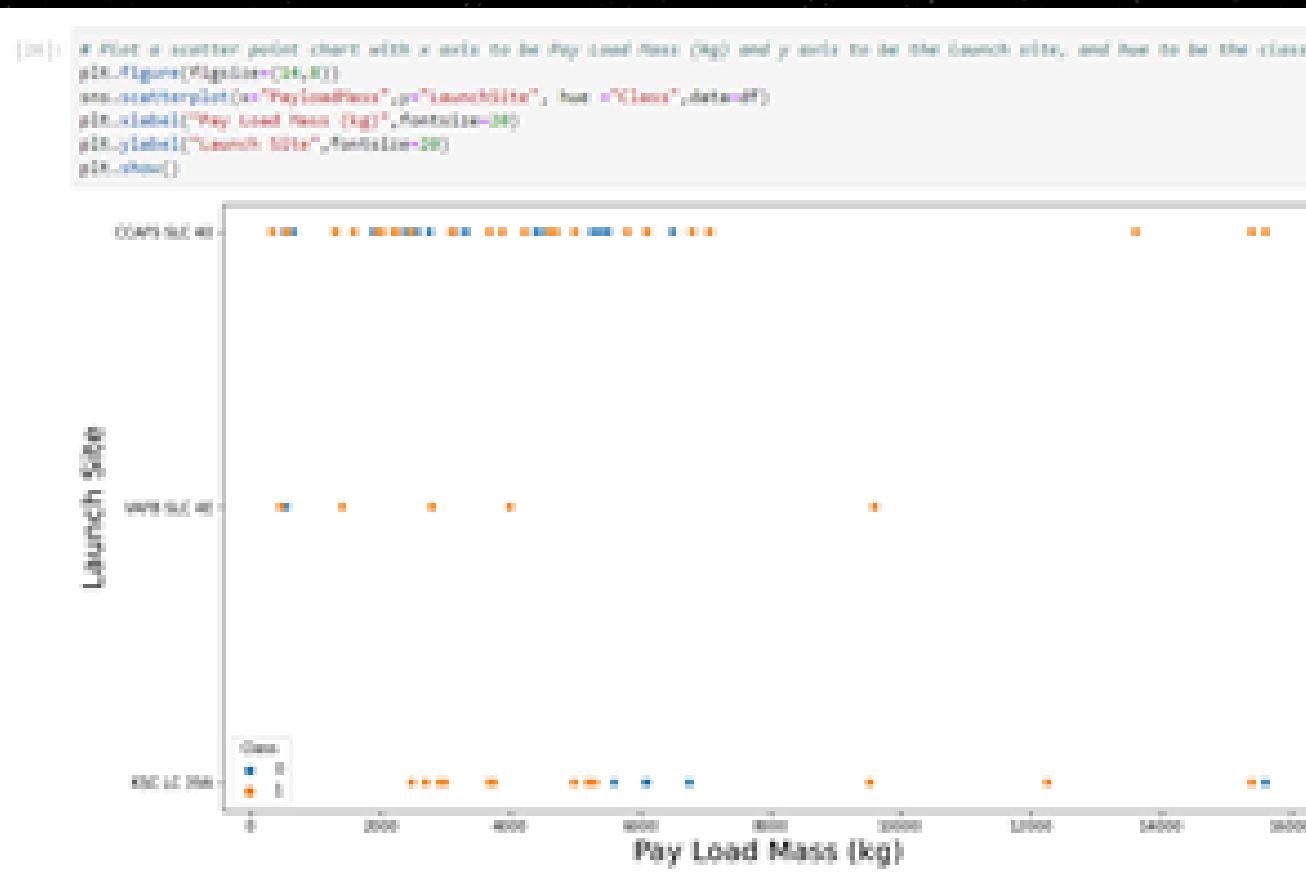
| Flight No. | Launch site | Payload | Payload mass | Orbit | Customer | Launch outcome | Version Booster | Booster landing | Date | Time |
|------------|-------------|--|--------------|-------|-----------|----------------|-----------------|-----------------|-----------------|-------|
| 0 | 1 | CCAFS Dragon Spacecraft Qualification Unit | 0 | LEO | SpaceX | Success\n | F9 v1.0B0003.1 | Failure | 4 June 2010 | 18:45 |
| 1 | 2 | CCAFS Dragon | 0 | LEO | NASA | Success | F9 v1.0B0004.1 | Failure | 8 December 2010 | 15:43 |
| 2 | 3 | CCAFS Dragon | 525 kg | LEO | NASA | Success | F9 v1.0B0005.1 | No attempt\n | 22 May 2012 | 07:44 |
| 3 | 4 | CCAFS SpaceX CRS-1 | 4,700 kg | LEO | NASA | Success\n | F9 v1.0B0006.1 | No attempt | 8 October 2012 | 00:35 |
| 4 | 5 | CCAFS SpaceX CRS-2 | 4,877 kg | LEO | NASA | Success\n | F9 v1.0B0007.1 | No attempt\n | 1 March 2013 | 15:10 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 116 | 117 | CCSFS Starlink | 15,600 kg | LEO | SpaceX | Success\n | F9 B5B1051.10 | Success | 9 May 2021 | 06:42 |
| 117 | 118 | KSC Starlink | ~14,000 kg | LEO | SpaceX | Success\n | F9 B5B1058.8 | Success | 15 May 2021 | 22:56 |
| 118 | 119 | CCSFS Starlink | 15,600 kg | LEO | SpaceX | Success\n | F9 B5B1063.2 | Success | 26 May 2021 | 18:59 |
| 119 | 120 | KSC SpaceX CRS-22 | 3,328 kg | LEO | NASA | Success\n | F9 B5B1067.1 | Success | 3 June 2021 | 17:29 |
| 120 | 121 | CCSFS SXM-8 | 7,000 kg | GTO | Sirius XM | Success\n | F9 B5 | Success | 6 June 2021 | 04:26 |

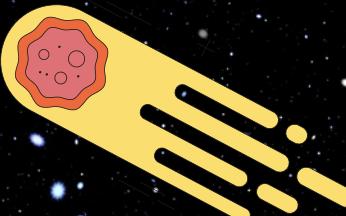




Data visualization: Exploratory

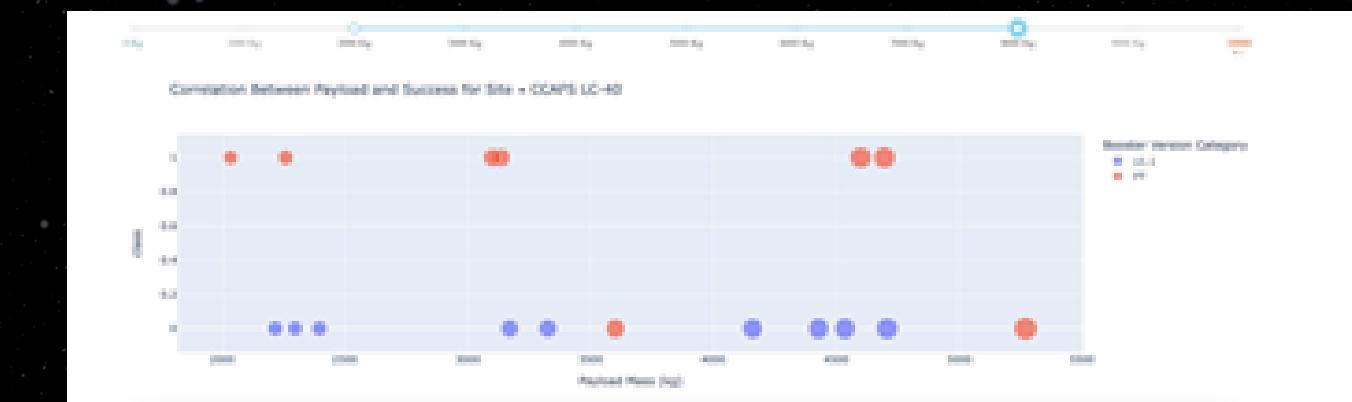
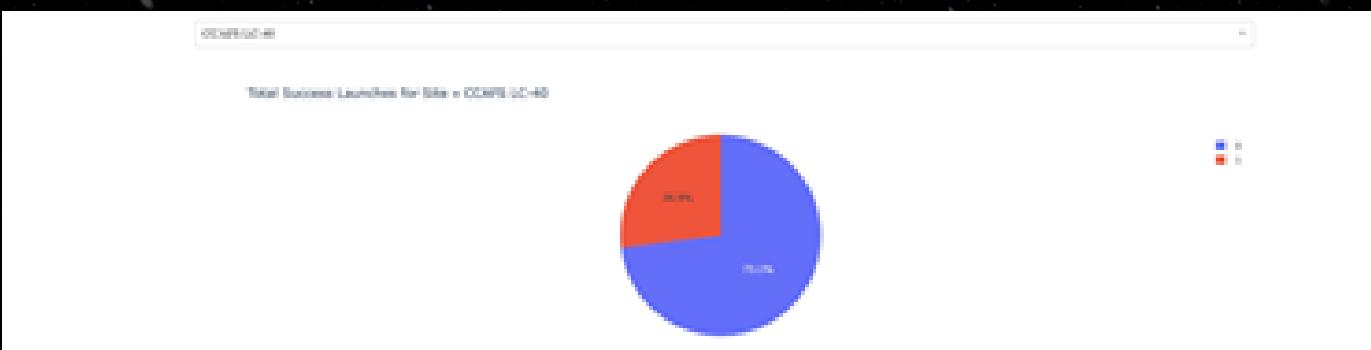
1. To understand relationships between variables of interest, using matplotlib and seaborn, collected and wrangled data was visualized using bar charts, line graphs and scatterplots.
2. To visualize and to mark launch sites where successful/ unsuccessful launches took place Folium was used to create interactive maps.

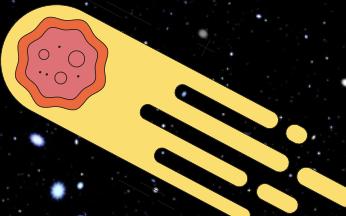




Data visualization: Dash

1. Using dash, an interactive pie and scatterplot chart was created to visualize the total success launches from each site and the correlation between payloadmass and mission outcome.





Machine learning

1. Using functions from scikit-learn, machine learning models were built. This was achieved by 1) standardizing the data 2) splitting the data into a train and test set 3) fitting built models on the training set 4) evaluating model accuracy based on the confusion matrix.

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
```

```
[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 10 test samples.

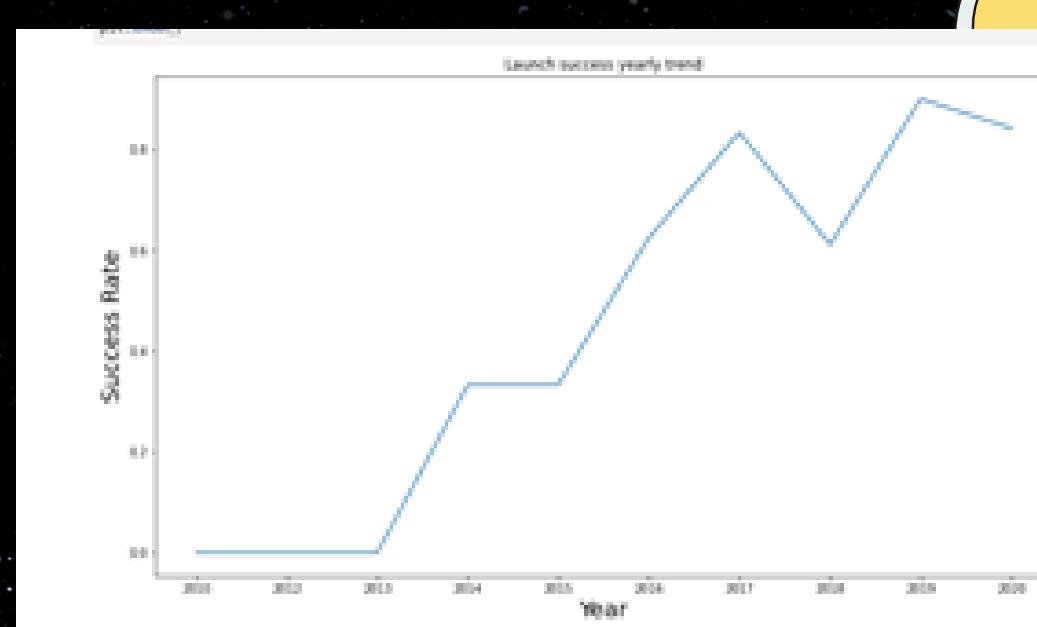
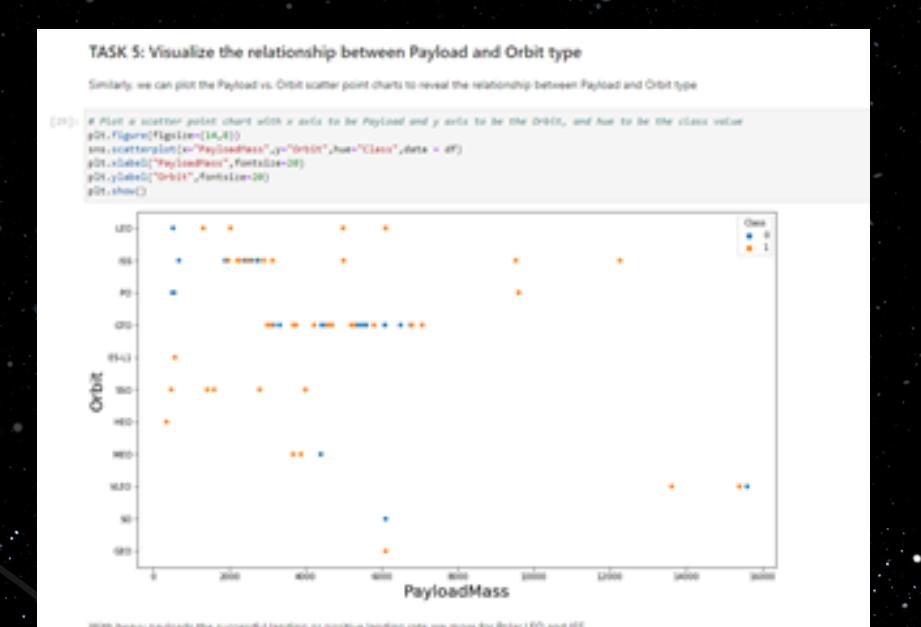
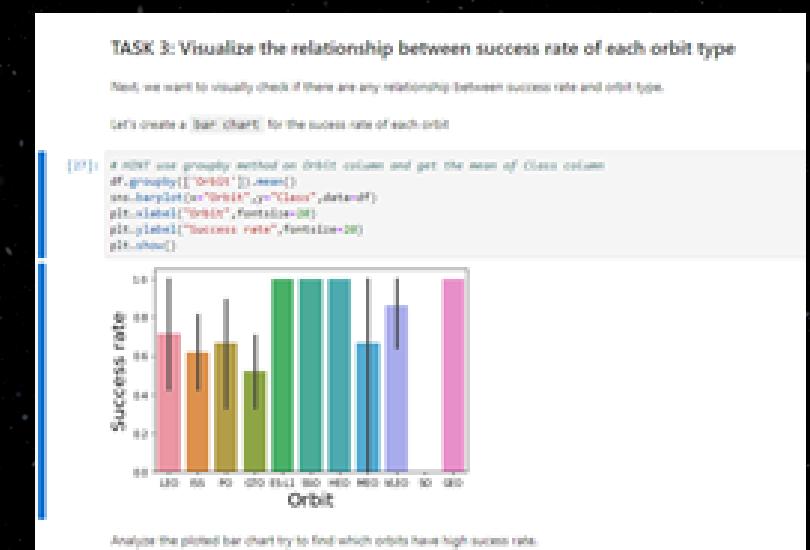
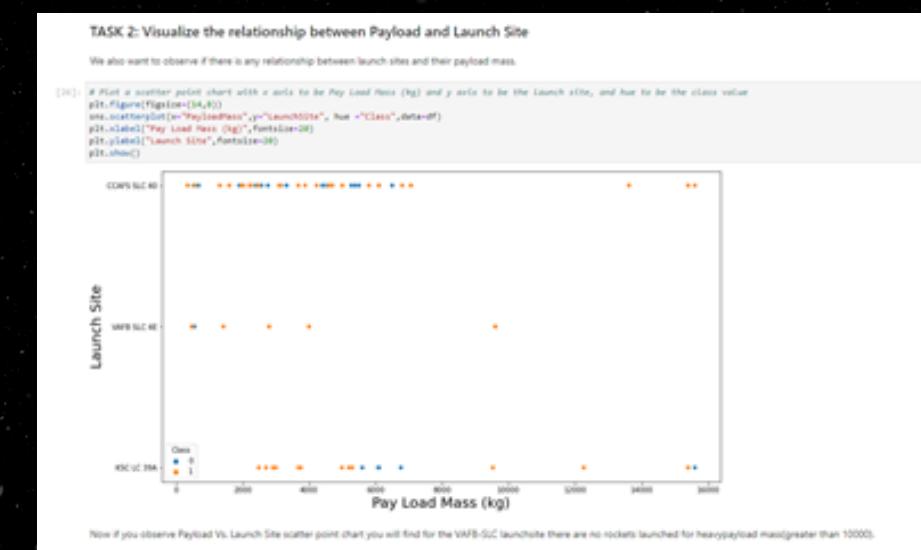
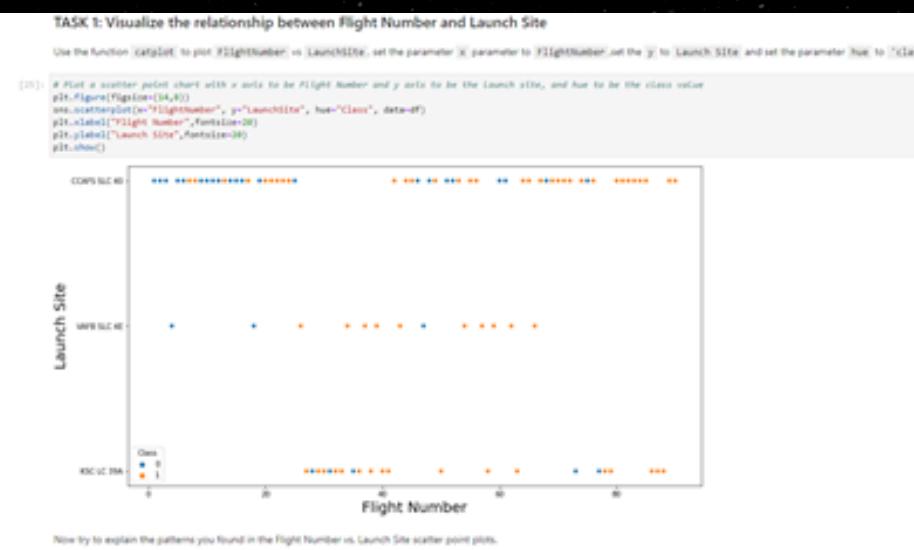
```
[1]: X_test.shape
```

```
[1]: (10,)
```



Results

EDA with visualization



Results

EDA with visualization

TASK 7: Create dummy variables to categorical columns

Use the function `get_dummies` and `Features_dataframe` to apply OneHotEncoder to the columns `Orbits`, `LaunchSite`, `LandingPad`, and `Serial`. Assign the value to the variable `Features_one_hot`, display the results using the method head. Your result `dataframe` must include all features including the encoded ones.

```

# apply the get_dummies() function on the categorical columns
features_one_hot = features

features_one_hot = pd.concat([features_one_hot,pd.get_dummies(df['Year'])], axis=1)
features_one_hot.drop(['Year'], axis=1, inplace=True)

features_one_hot = pd.concat([features_one_hot,pd.get_dummies(df['Unemployment'])], axis=1)
features_one_hot.drop(['Unemployment'], axis=1, inplace=True)

features_one_hot = pd.concat([features_one_hot,pd.get_dummies(df['Unemployed'])], axis=1)
features_one_hot.drop(['Unemployed'], axis=1, inplace=True)

features_one_hot = pd.concat([features_one_hot,pd.get_dummies(df['Vertical'])], axis=1)
features_one_hot.drop(['Vertical'], axis=1, inplace=True)

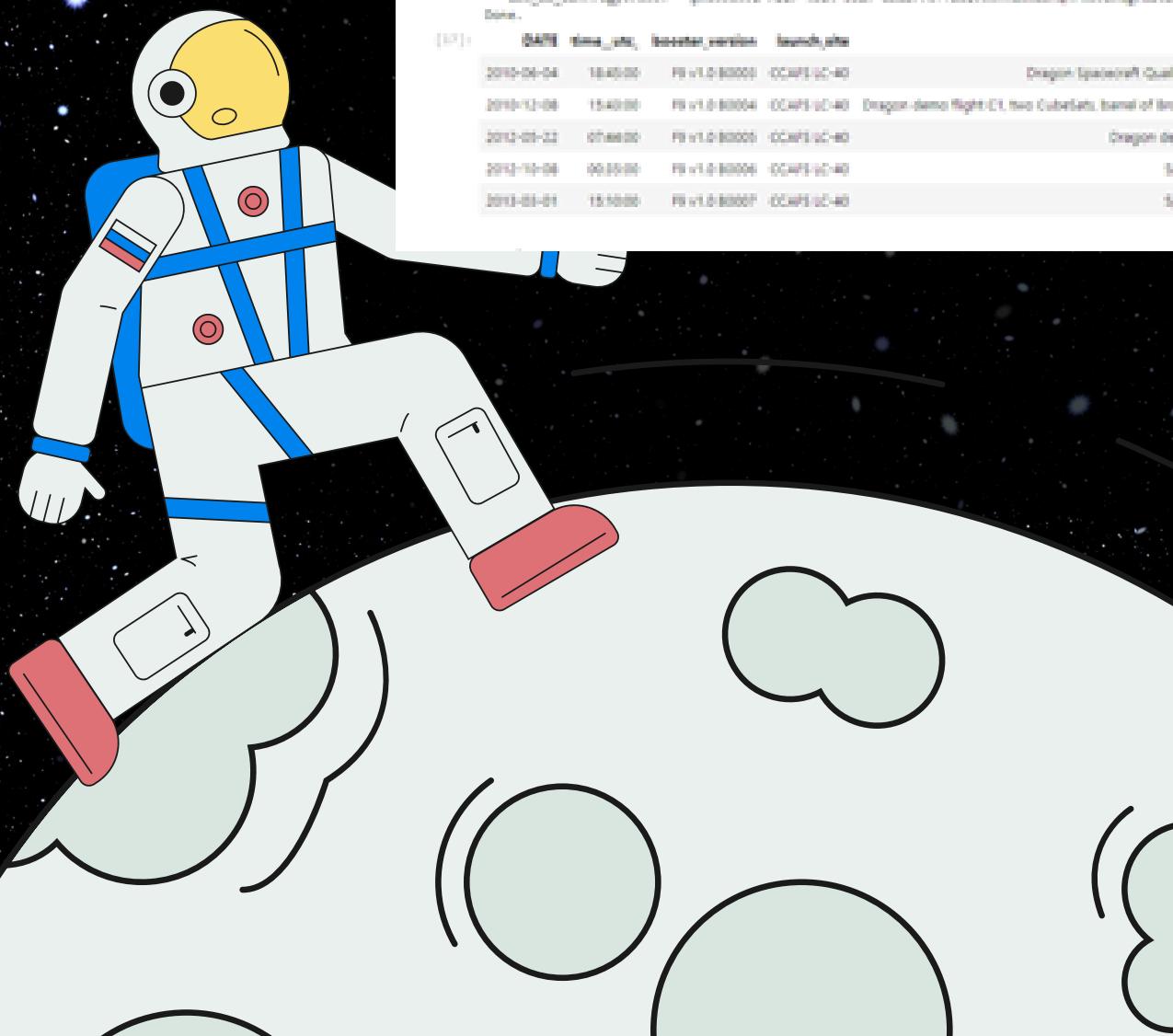
features_one_hot.drop(['Year'], axis=1, inplace=True)

```

1000 - 1000

Results

EDA with SQL



Task 1
Display the names of the unique launch sites in the space mission

```
[101] |sql| SELECT DISTINCT LAUNCH_SITE AS "Launch_Sites" FROM SPACEXFL
```

* Use_dbo..sys.tables::[dbo].[missions]::[missions]@0x0000000000000000..[dbo].[missions]@0x0000000000000000
Data

[102] |Launch_Sites|

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

Task 2
Display 5 records where launch sites begin with the string 'VCA'

```
[103] |sql| SELECT * FROM SPACEXFL WHERE LAUNCH_SITE LIKE "VCA%" LIMIT 5
```

* Use_dbo..sys.tables::[dbo].[missions]::[missions]@0x0000000000000000..[dbo].[missions]@0x0000000000000000
Data

[104] |

| DATE | TIME | LAUNCH_SITE | BOOSTER_VERSION | LAUNCH_POD | PAYLOAD | PAYLOAD_MASS_KG | ORBIT | CUSTOMER | MISSION_OUTCOME | LANDING_STATUS |
|------------|----------|---------------|-----------------|--|---------|-----------------|-----------------|----------|---------------------|----------------|
| 2019-09-04 | 10:40:00 | PF v1.0 80000 | CCAFS LC-40 | Dragon Spacecraft Qualification Unit | 0 | LEO | SpaceX | Success | Failure (parachute) | |
| 2019-12-08 | 15:40:00 | PF v1.0 80004 | CCAFS LC-40 | Dragon Demo Right C1, two Cubesats, band of brown cheese | 0 | LEO 650 | NASA (COTS) NAO | Success | Failure (parachute) | |
| 2020-01-22 | 07:40:00 | PF v1.0 80005 | CCAFS LC-40 | Oregon demo flight C2 | 120 | LEO 650 | NASA (COTS) | Success | No attempt | |
| 2020-01-26 | 09:30:00 | PF v1.0 80004 | CCAFS LC-40 | SpaceX DM-1 | 100 | LEO 650 | NASA (COTS) | Success | No attempt | |
| 2020-03-01 | 15:40:00 | PF v1.0 80007 | CCAFS LC-40 | SpaceX DM-2 | 477 | LEO 650 | NASA (COTS) | Success | No attempt | |

|

Task 3
Display the total payload mass carried by boosters launched by NASA (COTS)

```
[105] |sql| SELECT sum(payload_mass_kg) AS "Total payload mass by nasa (cots)" FROM SPACEXFL WHERE CUSTOMER = "NASA (COTS)"
```

* Use_dbo..sys.tables::[dbo].[missions]::[missions]@0x0000000000000000..[dbo].[missions]@0x0000000000000000
Data

[106] |Total payload mass by NASA (COTS)|

| |
|-------|
| 45000 |
|-------|

Task 4
Display average payload mass carried by booster version PF v1.0

```
[107] |sql| SELECT avg(payload_mass_kg) AS "Average payload mass by Booster Version PF v1.0" FROM SPACEXFL WHERE BOOSTER_VERSION = "PF v1.0"
```

* Use_dbo..sys.tables::[dbo].[missions]::[missions]@0x0000000000000000..[dbo].[missions]@0x0000000000000000
Data

[108] |Average payload mass by Booster Version PF v1.0|

| |
|------|
| 2020 |
|------|

Task 5
Get the date when the first successful landing outcome in ground pad was achieved.
Hint: Use min function

```
[109] |sql| SELECT min(DATE) AS "Date of first successful landing outcome in ground pad" FROM SPACEXFL WHERE MISSION_OUTCOME = "Success (ground pad)"
```

* Use_dbo..sys.tables::[dbo].[missions]::[missions]@0x0000000000000000..[dbo].[missions]@0x0000000000000000
Data

[110] |Date of first successful landing outcome in ground pad|

| |
|------------|
| 2019-12-23 |
|------------|

Task 6
Get the names of the boosters which have success in drone ship and have payload mass greater than 4000 but less than 6000

```
[111] |sql| SELECT BOOSTER_VERSION FROM SPACEXFL WHERE MISSION_OUTCOME = "Success (drone ship)" AND PAYLOAD_MASS_KG_BETWEEN 4000 AND 6000
```

* Use_dbo..sys.tables::[dbo].[missions]::[missions]@0x0000000000000000..[dbo].[missions]@0x0000000000000000
Data

[112] |Booster_version|

| |
|-------------|
| PF FT B1022 |
| PF FT B1028 |
| PF FT B1020 |
| PF FT B1002 |

Task 7
Get the total number of successful and failure mission outcomes

```
[113] |sql| SELECT number_of_success_outcomes, number_of_failure_outcomes FROM (SELECT COUNT(*) AS number_of_success_outcomes FROM SPACEXFL WHERE MISSION_OUTCOME LIKE "%success%") success_table, (SELECT COUNT(*) AS number_of_failure_outcomes FROM SPACEXFL WHERE MISSION_OUTCOME LIKE "%failure%") failure_table
```

* Use_dbo..sys.tables::[dbo].[missions]::[missions]@0x0000000000000000..[dbo].[missions]@0x0000000000000000
Data

[114] |number_of_success_outcomes| |number_of_failure_outcomes|

| | |
|-----|---|
| 100 | 1 |
|-----|---|

Task 8
Get the names of the booster_versions which have carried the maximum payload mass. Use a subquery

```
[115] |sql| SELECT DISTINCT BOOSTER_VERSION FROM SPACEXFL WHERE PAYLOAD_MASS_KG = (SELECT max(payload_mass_kg) FROM SPACEXFL)
```

* Use_dbo..sys.tables::[dbo].[missions]::[missions]@0x0000000000000000..[dbo].[missions]@0x0000000000000000
Data

[116] |Booster_version|

| |
|-------------|
| PF FT B1048 |
| PF FT B1048 |

Task 3
Display the total payload mass carried by boosters launched by NASA (COTS)

```
[105] |sql| SELECT sum(payload_mass_kg) AS "Total payload mass by nasa (cots)" FROM SPACEXFL WHERE CUSTOMER = "NASA (COTS)"
```

* Use_dbo..sys.tables::[dbo].[missions]::[missions]@0x0000000000000000..[dbo].[missions]@0x0000000000000000
Data

[106] |Total payload mass by NASA (COTS)|

| |
|-------|
| 45000 |
|-------|

Task 4
Display average payload mass carried by booster version PF v1.0

```
[107] |sql| SELECT avg(payload_mass_kg) AS "Average payload mass by Booster Version PF v1.0" FROM SPACEXFL WHERE BOOSTER_VERSION = "PF v1.0"
```

* Use_dbo..sys.tables::[dbo].[missions]::[missions]@0x0000000000000000..[dbo].[missions]@0x0000000000000000
Data

[108] |Average payload mass by Booster Version PF v1.0|

| |
|------|
| 2020 |
|------|

Task 5
Get the date when the first successful landing outcome in ground pad was achieved.
Hint: Use min function

```
[109] |sql| SELECT min(DATE) AS "Date of first successful landing outcome in ground pad" FROM SPACEXFL WHERE MISSION_OUTCOME = "Success (ground pad)"
```

* Use_dbo..sys.tables::[dbo].[missions]::[missions]@0x0000000000000000..[dbo].[missions]@0x0000000000000000
Data

[110] |Date of first successful landing outcome in ground pad|

| |
|------------|
| 2019-12-23 |
|------------|

Task 6
Get the names of the boosters which have success in drone ship and have payload mass greater than 4000 but less than 6000

```
[111] |sql| SELECT BOOSTER_VERSION FROM SPACEXFL WHERE MISSION_OUTCOME = "Success (drone ship)" AND PAYLOAD_MASS_KG_BETWEEN 4000 AND 6000
```

* Use_dbo..sys.tables::[dbo].[missions]::[missions]@0x0000000000000000..[dbo].[missions]@0x0000000000000000
Data

[112] |Booster_version|

| |
|-------------|
| PF FT B1022 |
| PF FT B1028 |
| PF FT B1020 |
| PF FT B1002 |

Task 7
Get the total number of successful and failure mission outcomes

```
[113] |sql| SELECT number_of_success_outcomes, number_of_failure_outcomes FROM (SELECT COUNT(*) AS number_of_success_outcomes FROM SPACEXFL WHERE MISSION_OUTCOME LIKE "%success%") success_table, (SELECT COUNT(*) AS number_of_failure_outcomes FROM SPACEXFL WHERE MISSION_OUTCOME LIKE "%failure%") failure_table
```

* Use_dbo..sys.tables::[dbo].[missions]::[missions]@0x0000000000000000..[dbo].[missions]@0x0000000000000000
Data

[114] |number_of_success_outcomes| |number_of_failure_outcomes|

| | |
|-----|---|
| 100 | 1 |
|-----|---|

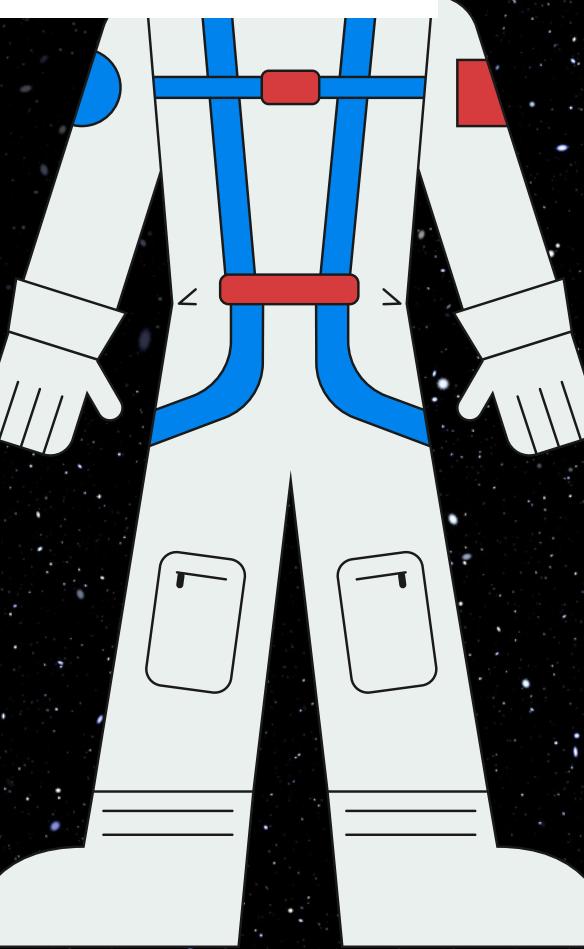
Task 8
Get the names of the booster_versions which have carried the maximum payload mass. Use a subquery

```
[115] |sql| SELECT DISTINCT BOOSTER_VERSION FROM SPACEXFL WHERE PAYLOAD_MASS_KG = (SELECT max(payload_mass_kg) FROM SPACEXFL)
```

* Use_dbo..sys.tables::[dbo].[missions]::[missions]@0x0000000000000000..[dbo].[missions]@0x0000000000000000
Data

[116] |Booster_version|

| |
|-------------|
| PF FT B1048 |
| PF FT B1048 |



Results

EDA with SQL

Task 9
List the failed landing_outcomes in drone ship, their booster versions, and launch site names for the year 2015

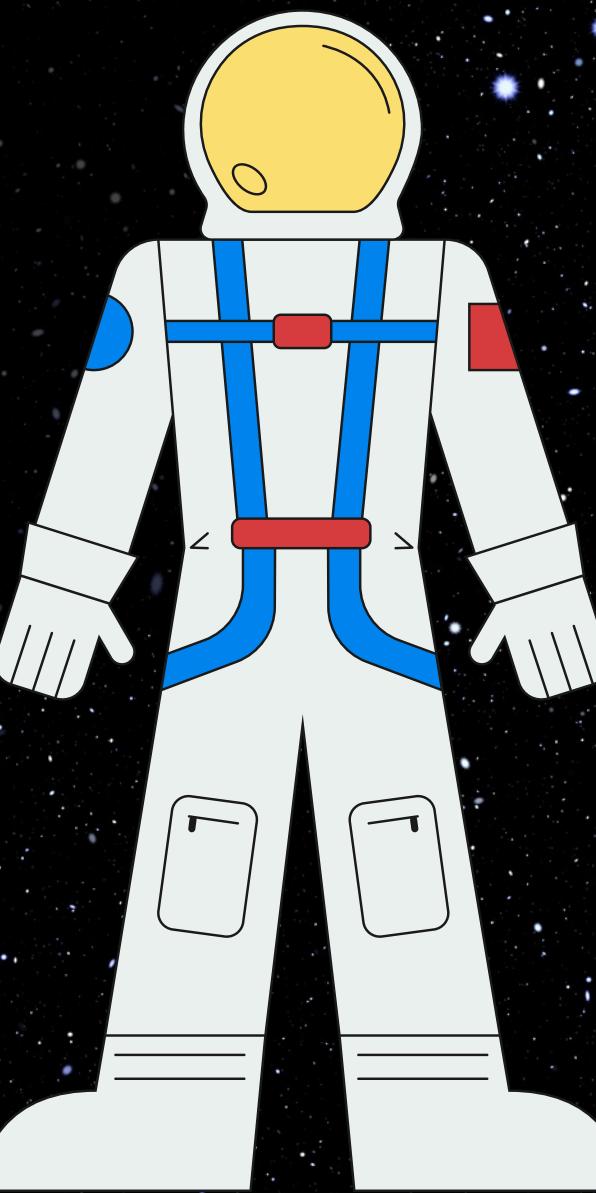
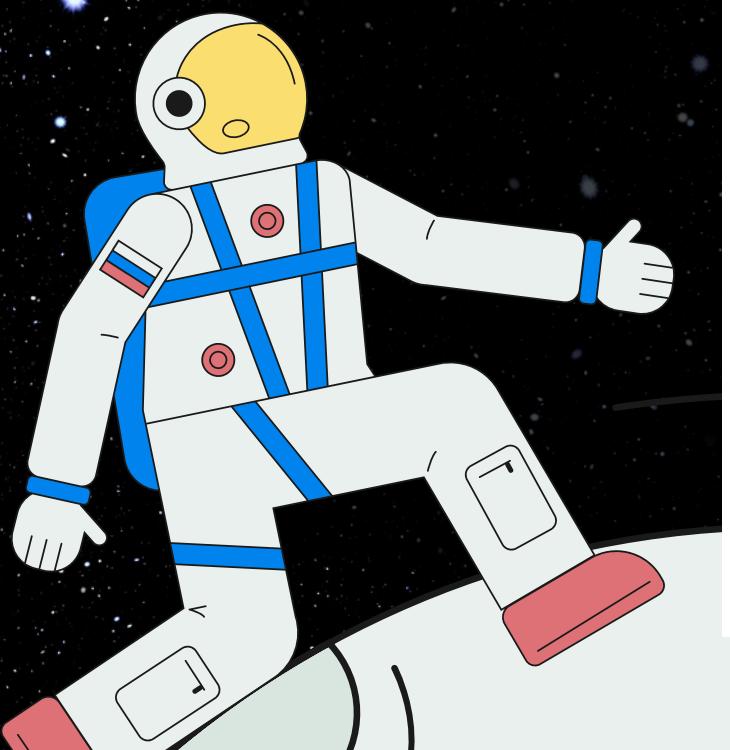
```
SQL SELECT DATE, BOOSTER_VERSION, LAUNCH_SITE FROM SPACEXTEL WHERE year(DATE) = '2015' AND landing_outcome = 'Failure (drone ship)';  
* See also https://github.com/justmarkham/PythonDataScienceHandbook/blob/master/notebooks/08_sql.ipynb
```

| DATE | booster_version | Launch_site |
|------------|-----------------|-------------|
| 2015-06-10 | FL-1.1 B9012 | CCAFS LC-40 |
| 2015-08-16 | FL-1.1 B9013 | CCAFS LC-40 |

Task 10
Rank the count of landing outcomes (such as Failure (drone ship) or Success (ground pad)) between the date 2010-01-01 and 2017-09-26, in descending order

```
SQL SELECT landing_outcome, COUNT(landing_outcome) AS Landing_Count FROM SPACEXTEL WHERE DATE BETWEEN '2010-01-01' AND '2017-09-26' GROUP BY landing_outcome ORDER BY COUNT(landing_outcome) DESC;  
* See also https://github.com/justmarkham/PythonDataScienceHandbook/blob/master/notebooks/08_sql.ipynb
```

| landing_outcome | landing_count |
|-------------------------|---------------|
| No attempt | 10 |
| Failure (drone ship) | 1 |
| Success (drone ship) | 10 |
| Controlled (ocast) | 1 |
| Success (ground pad) | 10 |
| Failure (personal) | 1 |
| Uncontrolled (ocast) | 1 |
| Prohibited (drone ship) | 1 |



Results

Interactive maps with folium

Task 1: Mark all launch sites on a map

First, let's try to add each site's location on a map using site's latitude and longitude coordinates.

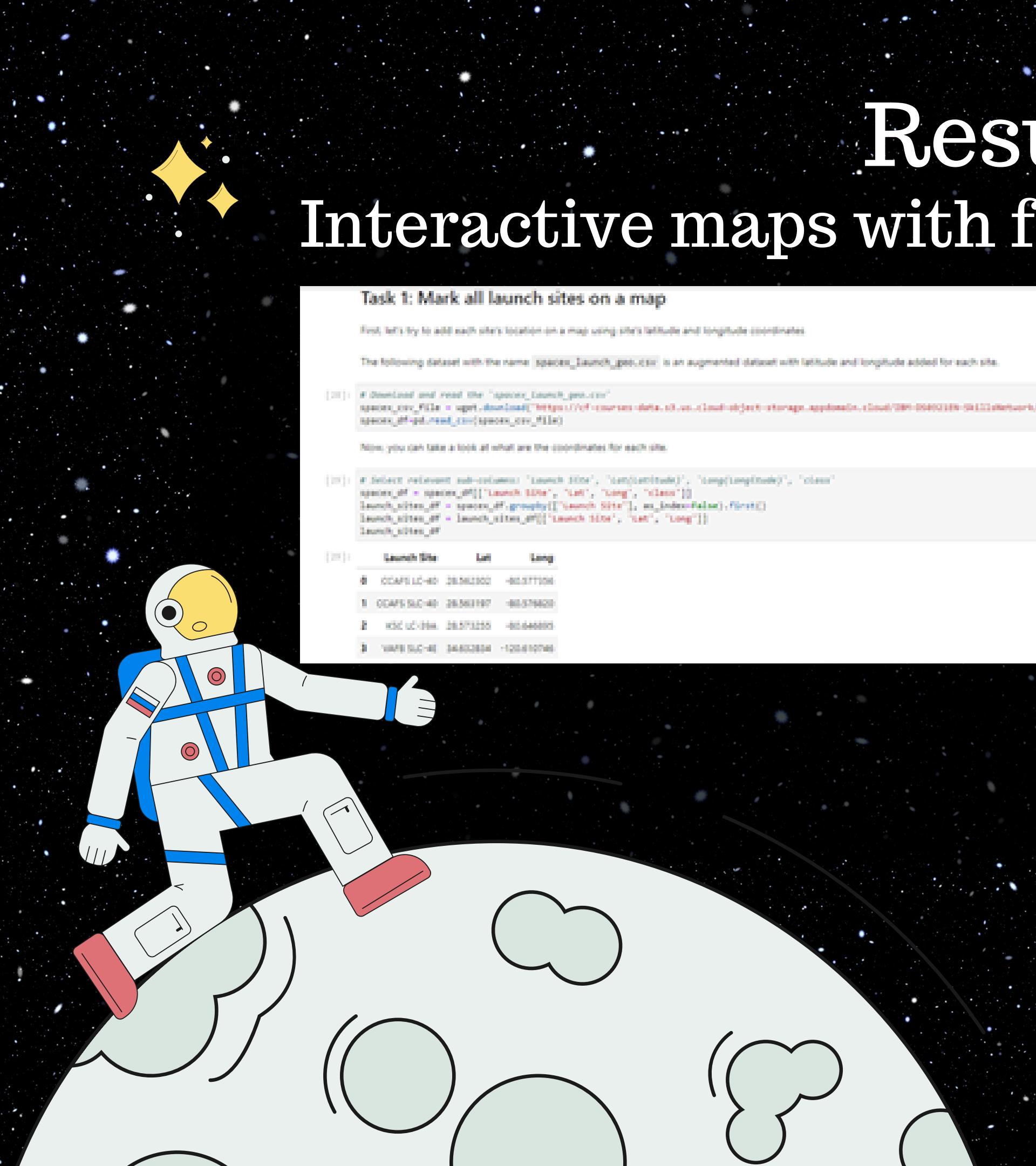
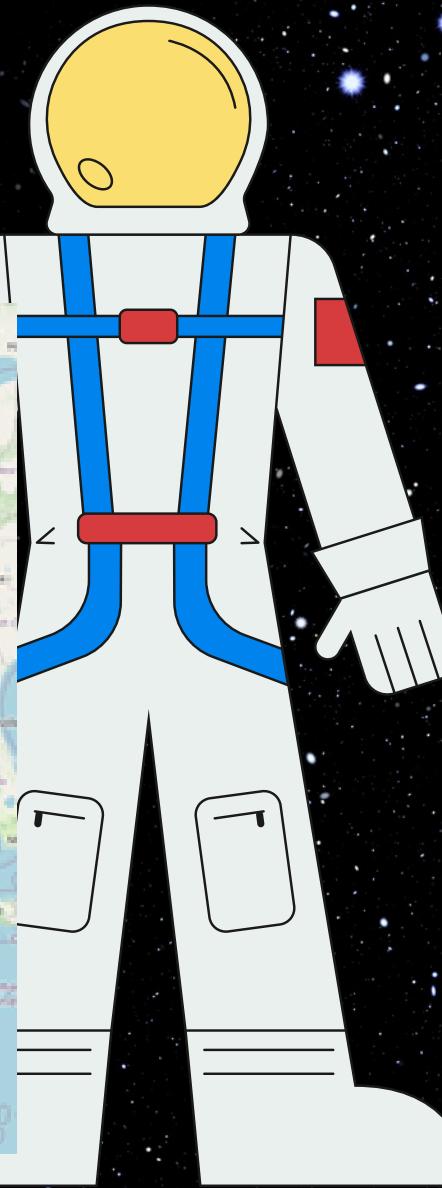
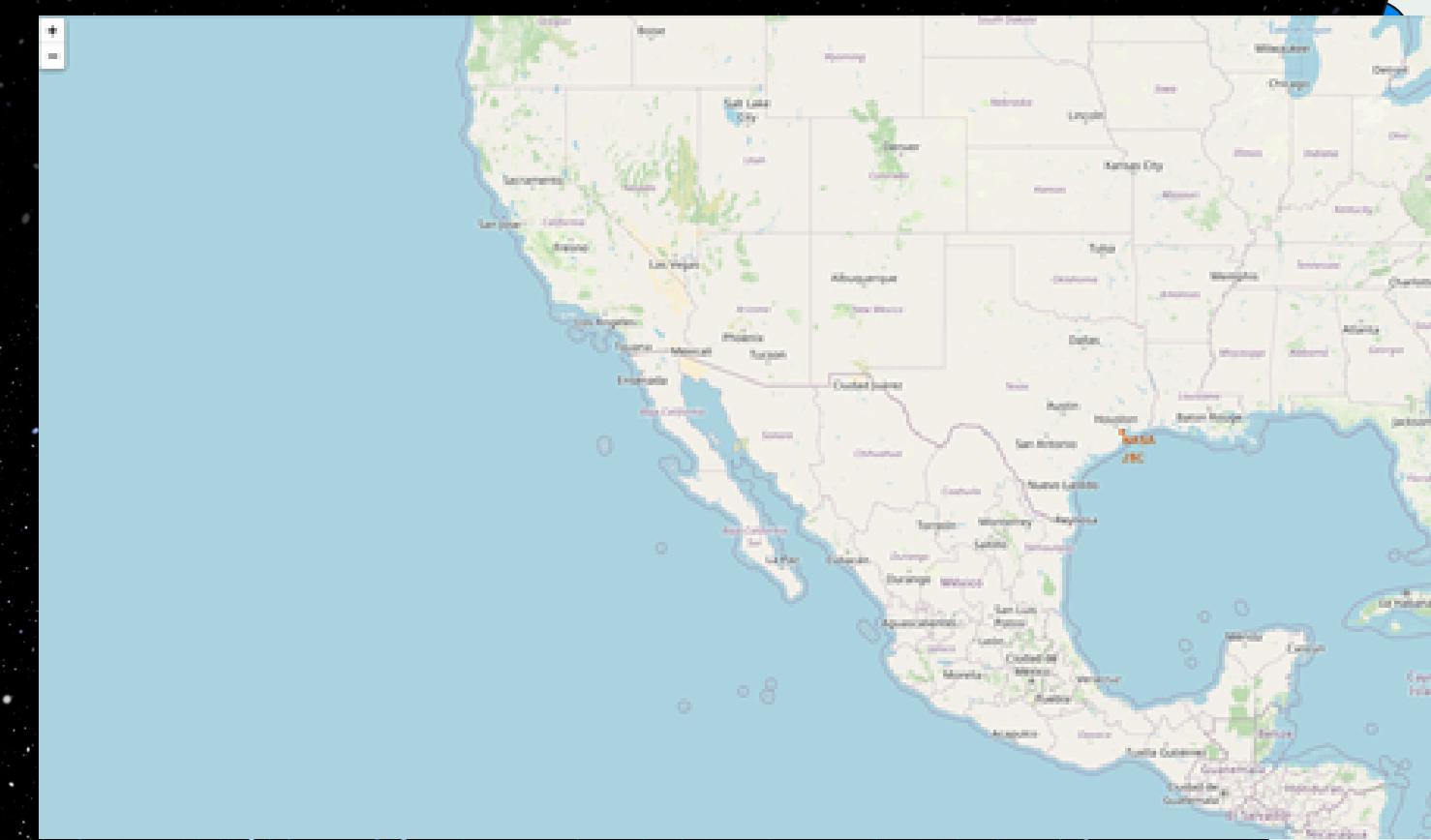
The following dataset with the name `spaceLaunch_gov.csv` is an augmented dataset with latitude and longitude added for each site.

```
[1]: # Download and read the "spaceLaunch_gov.csv"
space_gov_file = wget.download("https://cf-courses-data.s3.us.cloud-object-storage.applications/cf-000-004-008-0012/starter-data/spaceLaunch_gov.csv")
space_gov_df = pd.read_csv(space_gov_file)
```

Now, you can take a look at what are the coordinates for each site.

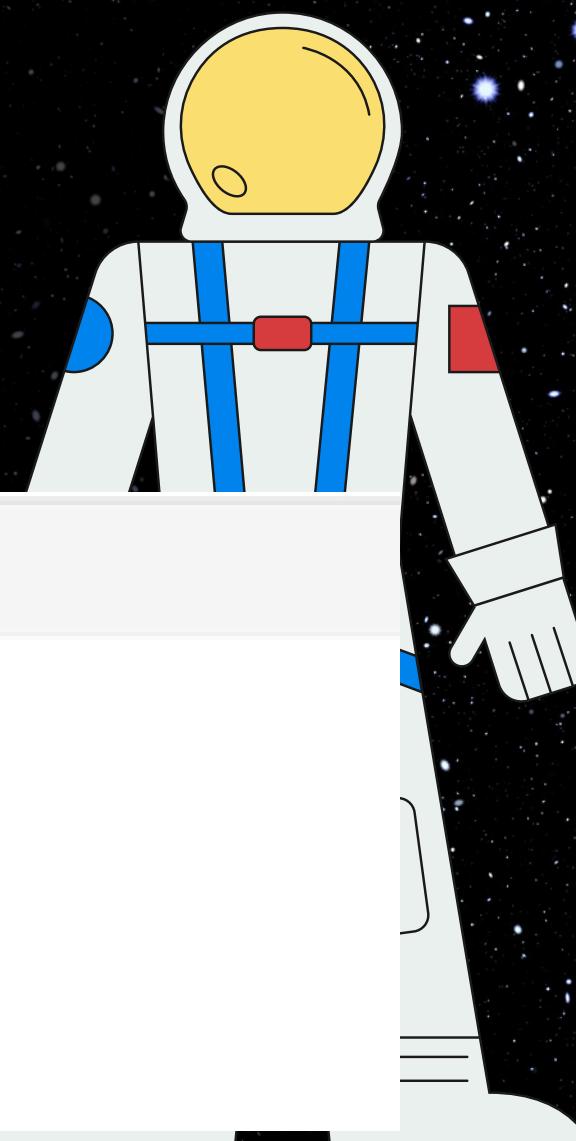
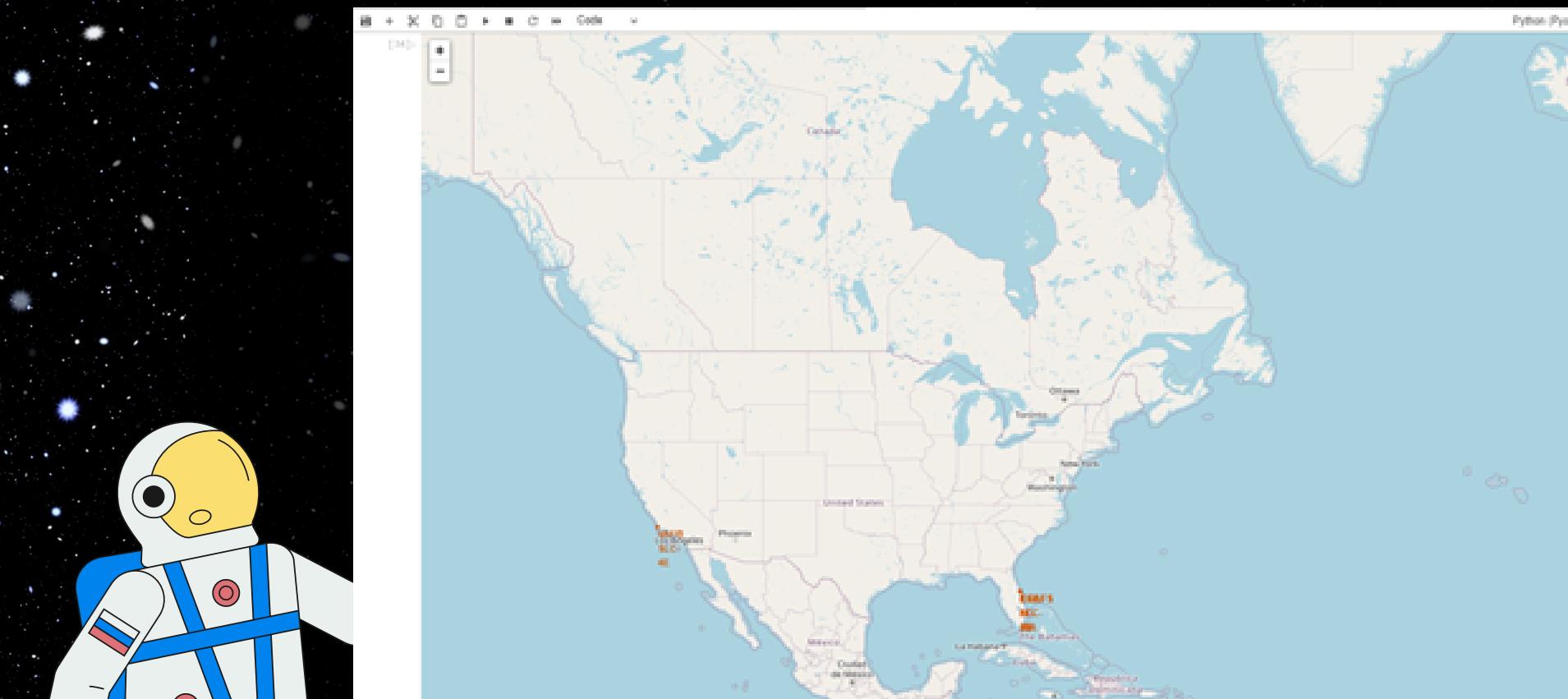
```
[2]: # Select relevant sub-columns: "Launch Site", "Lat", "Long"
space_gov_df[['Launch Site', 'Lat', 'Long']]
# Launch_Site      Lat      Long
# 0     CCAS LC-40  31.963852 -80.377500
# 1     CCAS SLC-40  38.631197 -80.779429
# 2     KSC LC-39A  28.571193 -80.444993
# 3     VAFB SLC-4  34.632084 -120.410793
```

| Launch Site | Lat | Long |
|-------------|-----------|-------------|
| CCAS LC-40 | 31.963852 | -80.377500 |
| CCAS SLC-40 | 38.631197 | -80.779429 |
| KSC LC-39A | 28.571193 | -80.444993 |
| VAFB SLC-4 | 34.632084 | -120.410793 |



Results

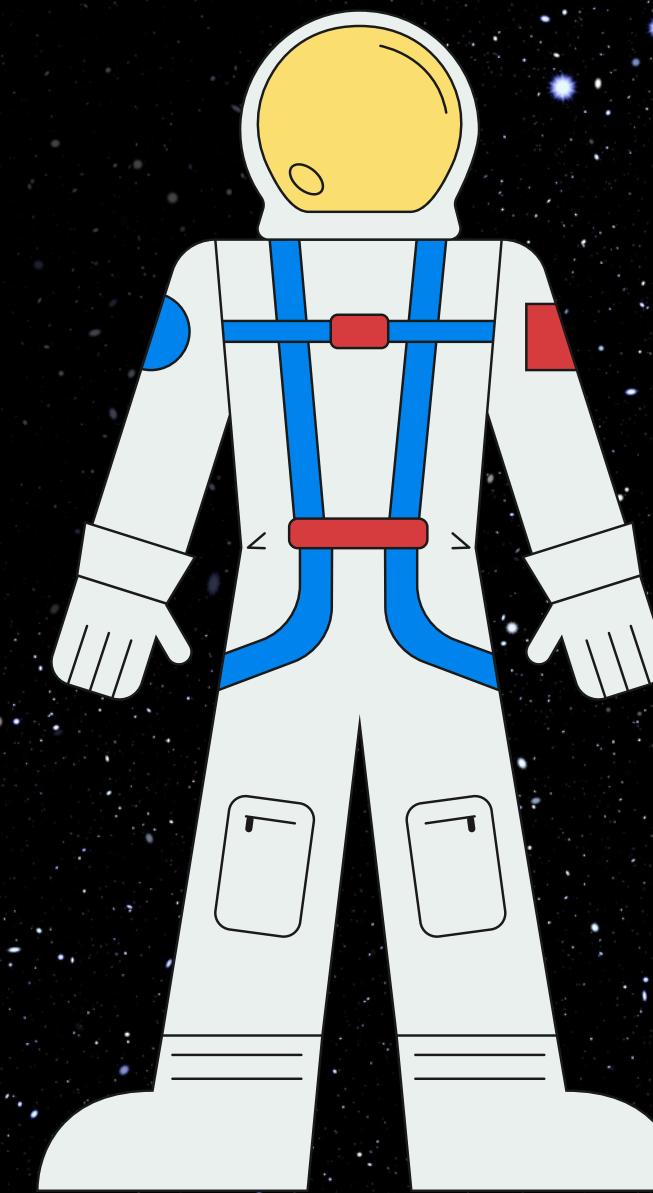
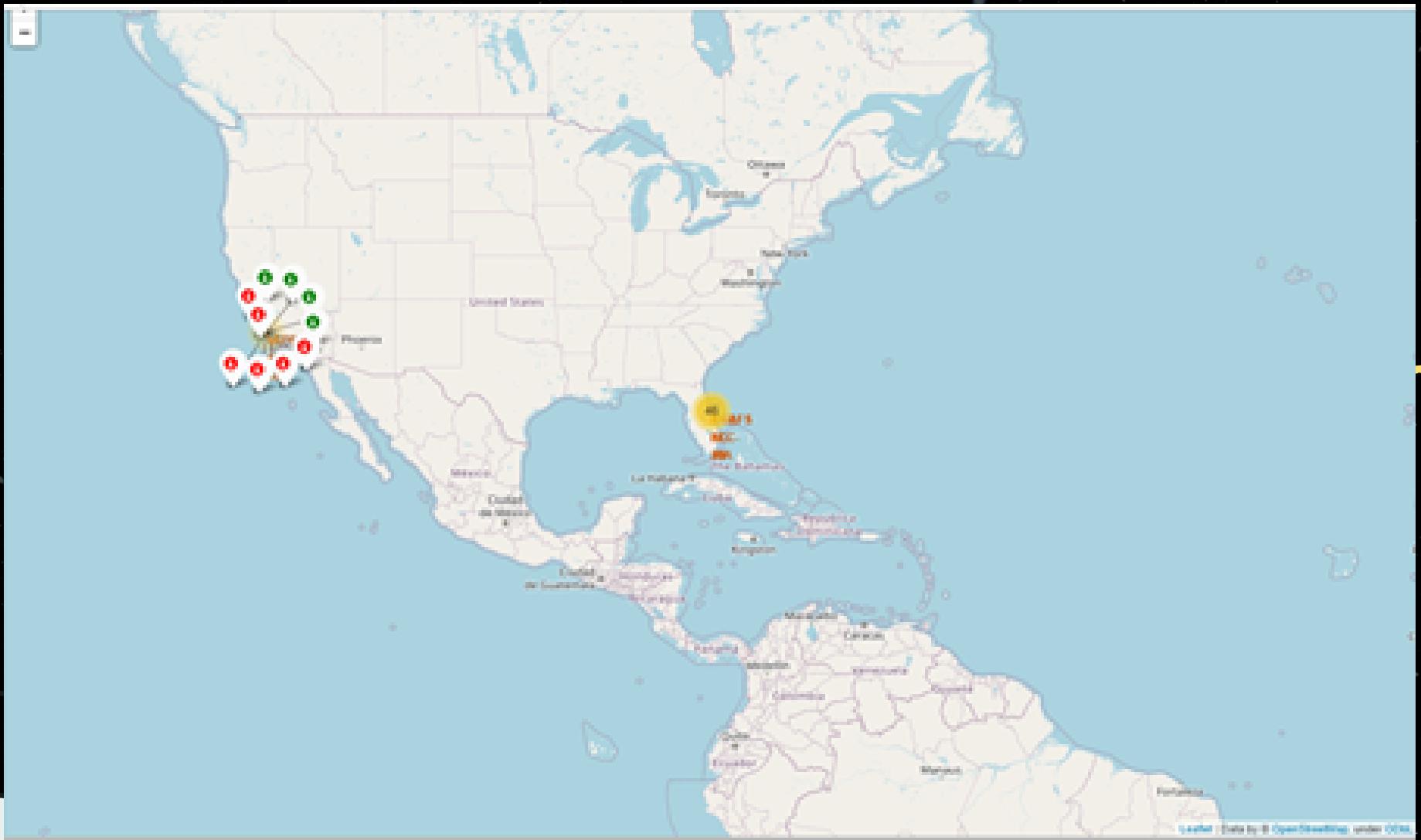
Interactive maps with folium



| # | Launch Site | Lat | Long | Class | marker_color |
|----|--------------|----------|-----------|-------|--------------|
| 46 | KSC LC-39A | 28.57021 | -80.57021 | 1 | green |
| 47 | KSC LC-39A | 28.57021 | -80.56999 | 1 | green |
| 48 | KSC LC-39A | 28.57021 | -80.56999 | 1 | green |
| 49 | CCAFS SLC-40 | 28.56997 | -80.57982 | 1 | green |
| 50 | CCAFS SLC-40 | 28.56997 | -80.57982 | 1 | green |
| 51 | CCAFS SLC-40 | 28.56997 | -80.57982 | 0 | red |
| 52 | CCAFS SLC-40 | 28.56997 | -80.57982 | 0 | red |
| 53 | CCAFS SLC-40 | 28.56997 | -80.57982 | 0 | red |
| 54 | CCAFS SLC-40 | 28.56997 | -80.57982 | 1 | green |
| 55 | CCAFS SLC-40 | 28.56997 | -80.57982 | 0 | red |

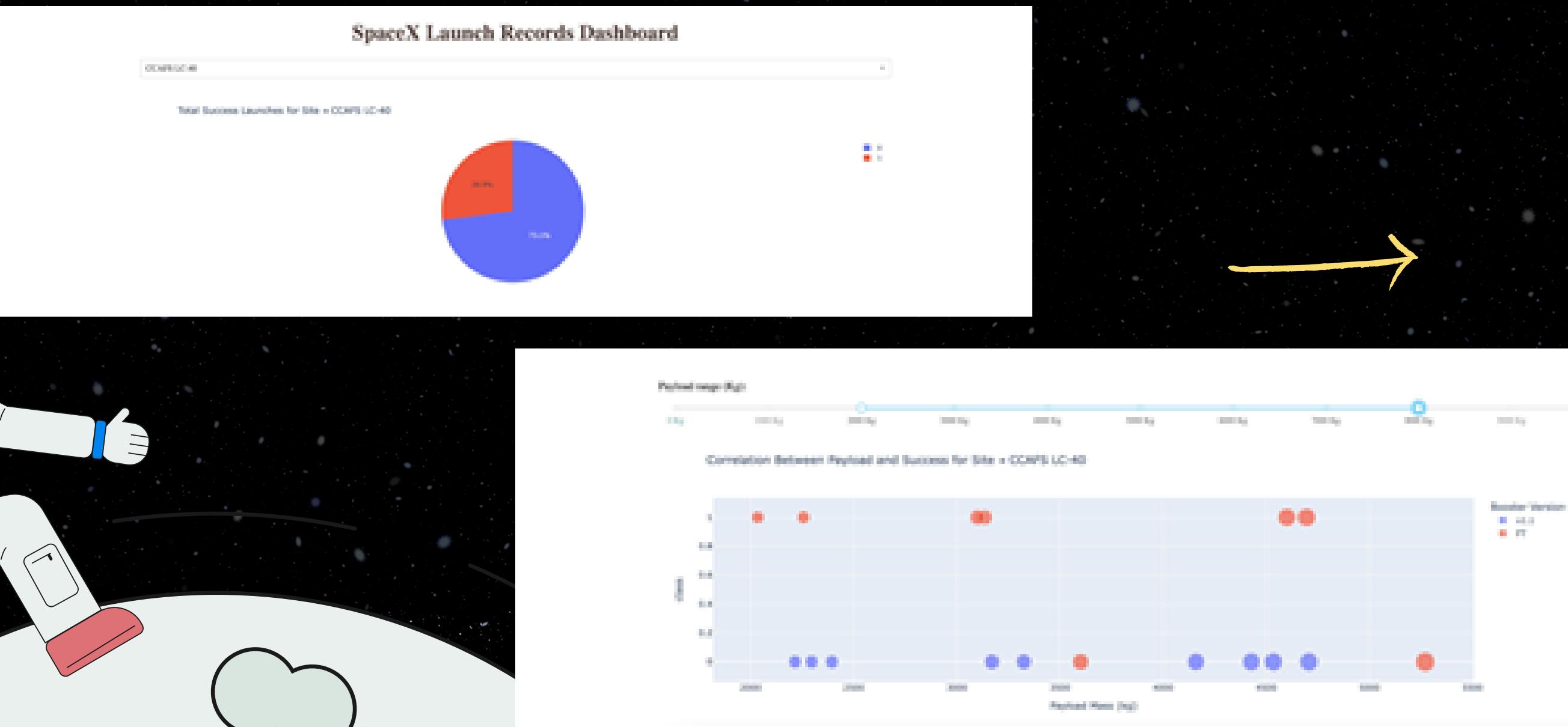
Results

Interactive maps with folium



Results

Dashboard



Results

Predictive analyses

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].

```
[1]: parameters = [{"C": 0.01, "penalty": "l1"},  
                 {"C": 0.01, "penalty": "l2"},  
                 {"C": 0.01, "solver": "liblinear"},  
                 {"C": 0.01, "solver": "lbfgs"},  
                 {"C": 0.01, "solver": "newton-cg"},  
                 {"C": 0.01, "solver": "sag"},  
                 {"C": 0.01, "solver": "saga"}]  
  
[2]: logreg = LogisticRegression()  
logreg_cv = GridSearchCV(logreg, parameters, scoring="accuracy", cv=10)  
logreg_cv.fit(X_train, Y_train)
```

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_].

```
[3]: print("Best hyperparameters : (best parameters) : ", logreg_cv.best_params_)  
print("Accuracy : ", logreg_cv.best_score_)  
  
best hyperparameters : (best parameters) : {'C': 0.01, 'penalty': 'l2', 'solver': 'lbfgs'}  
accuracy : 0.8888888888888889
```



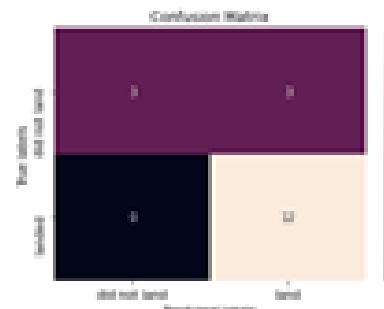
TASK 5

Calculate the accuracy on the test data using the method [score].

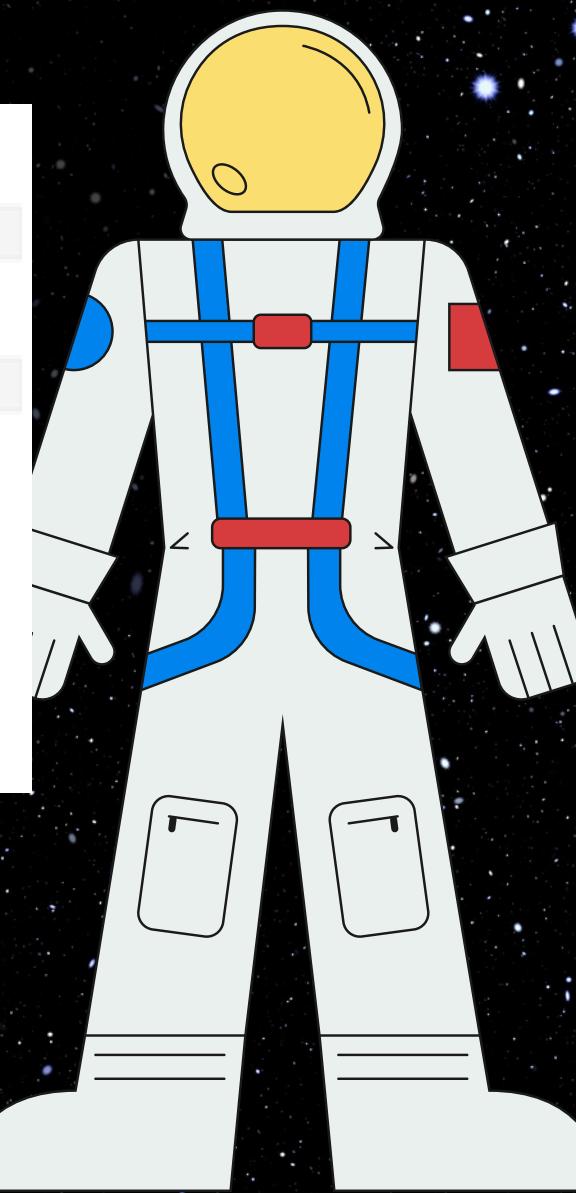
```
[1]: X_test_accuracy = logreg_cv.score(X_test, Y_test)  
X_test_accuracy
```

Let's look at the confusion matrix.

```
[2]: y_hat = logreg_cv.predict(X_test)  
plot_confusion_matrix(Y_test, y_hat)
```



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



Results

Predictive analyses

TASK 6

Create a support vector machine object then create a `GridSearchCV` object `cv_grid` with `cv=5`. Use the object to find the best parameters from the dictionary `parameters`.

```
[1]: parameters = [{"kernel": "linear", "C": 1}, {"kernel": "poly", "C": 1, "gamma": "scale"}, {"kernel": "rbf", "C": 1, "gamma": "scale"}, {"kernel": "sigmoid", "C": 1, "gamma": "scale"}]
cv_grid = GridSearchCV()

[2]: cv_grid = GridSearchCV(cv=cv_grid, parameters, scoring='accuracy', cv=5)
cv_grid.fit(X_train, Y_train)

[3]: print("Final Hyperparameters : (best parameters): ", cv_grid.best_params_)
print("Accuracy : ", cv_grid.best_score_)

Final Hyperparameters : (best parameters): {'C': 1.0, 'gamma': 'scale', 'kernel': 'rbf'}
accuracy : 0.8994999999999999
```

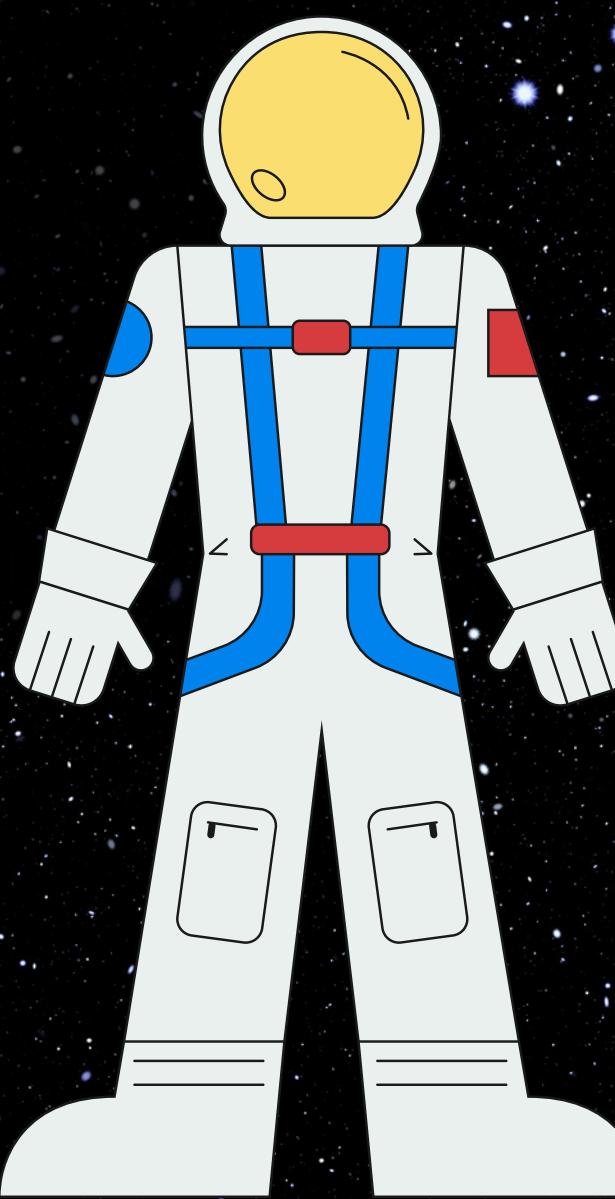
TASK 7

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

```
[24]: acc_accuracy = cv_grid.score(X_test, Y_test)
acc_accuracy
```

We can plot the confusion matrix:

```
[25]: plot_confusion_matrix(X_test)
plot_confusion_matrix(Y_test, y_hat)
```



Results

Predictive analyses

TASK 8

Create a decision tree classifier object then create a `DecisionTreeClassifier` object `tree_cv` with $n = 10$. Ask the object to find the best parameters (use `accuracy` parameter).

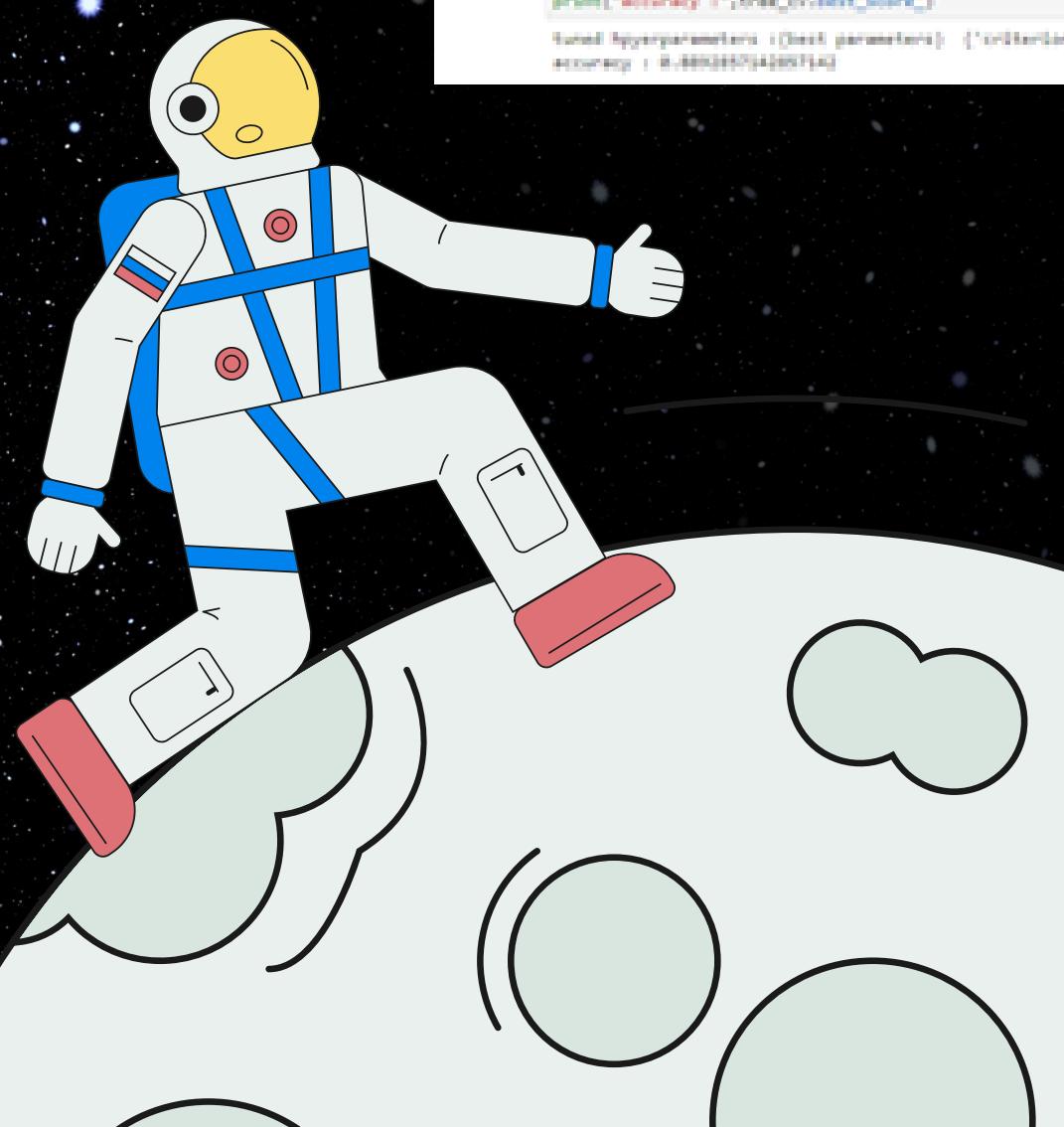
```
[1]: parameters = {"criterion": ["gini", "entropy"],
                 "splitter": ["best", "random"],
                 "max_depth": [2, 4, 6, 8, 10, 12],
                 "min_features": [2, 4, 6, 8],
                 "min_samples_leaf": [1, 2, 4],
                 "min_samples_split": [2, 4, 6]}

tree = DecisionTreeClassifier()

[2]: tree_cv = GridSearchCV(tree, parameters, scoring="accuracy", cv=10)
tree_cv = tree_cv.fit(X_train, Y_train)

[3]: print("Best Hyperparameters (test parameters):", tree_cv.best_params_)
print("Accuracy:", tree_cv.best_score_)

Best Hyperparameters (test parameters): {'criterion': 'gini', 'max_depth': 4, 'min_features': 'sqrt', 'min_samples_leaf': 1, 'min_samples_split': 2, 'splitter': 'best'}
accuracy: 0.8333333333333333
```



TASK 9

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

```
[1]: tree_accuracy = tree_cv.score(X_test, Y_test)
tree_accuracy

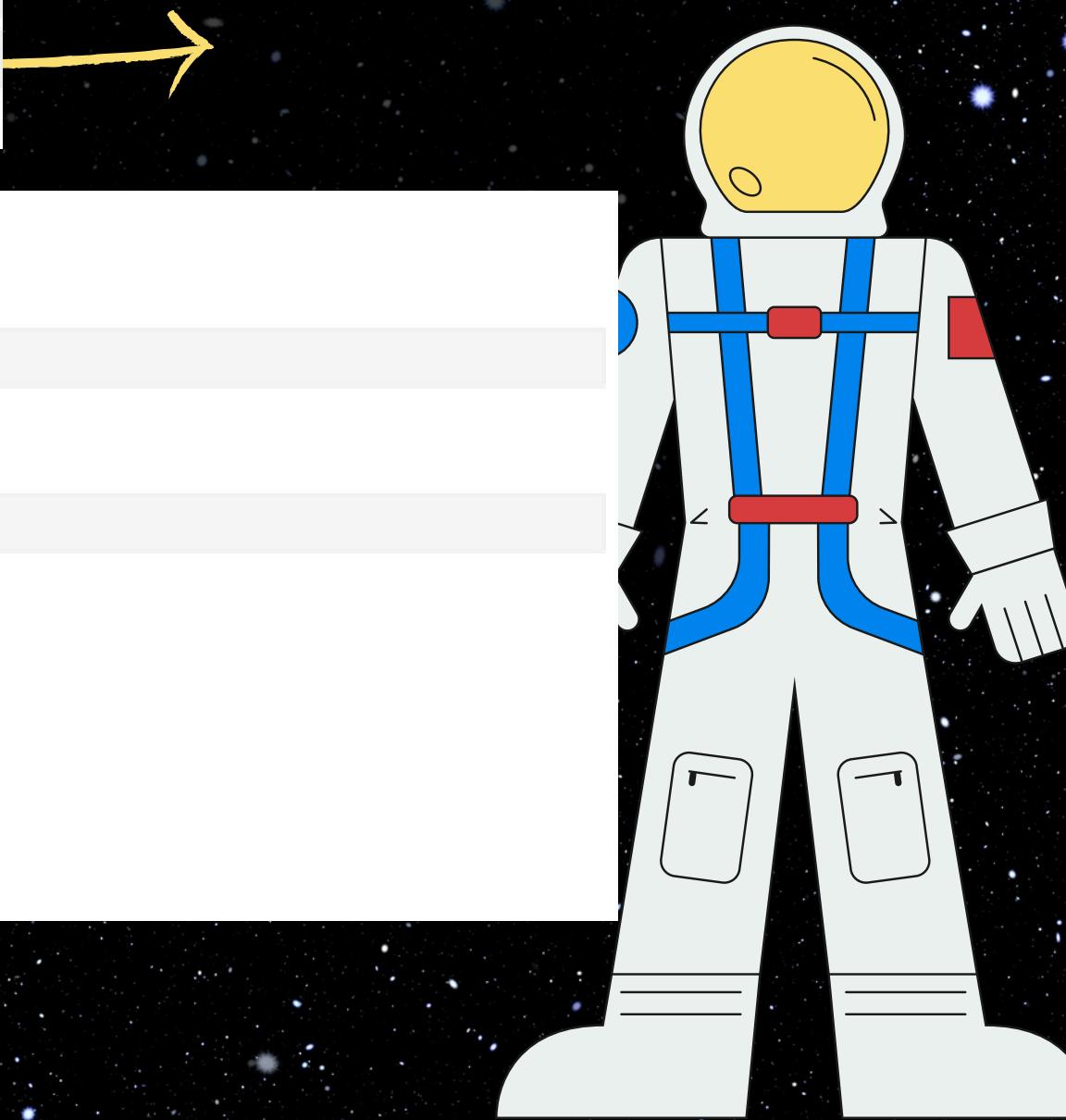
[2]: 0.8333333333333333
```

We can plot the confusion matrix.

```
[3]: y_hat = tree_cv.predict(X_test)
plot_confusion_matrix(Y_test, y_hat)
```

Confusion Matrix

| | | Predicted labels |
|-------------|---------------|------------------|
| True labels | Actual labels | |
| | 0 | 1 |
| 0 | 80 | 10 |
| 1 | 10 | 90 |



Results

Predictive analyses

TASK 10

Create a `k_nearest_neighbors` object then create a `GridSearchCV` object `base_cv` with $n=10$. Fit the object to find the best parameters from the dictionary `parameters`.

```
[1]: parameters = {"n_neighbors": [1, 2, 3, 4, 5, 6, 7, 8, 9, 10],  
                 "algorithm": ["auto", "ball_tree", "kd_tree", "brute"],  
                 "p": [1,2]}  
  
knn = KNeighborsClassifier()  
  
[2]: base_cv = GridSearchCV(knn, parameters, scoring='accuracy', cv=10)  
base_cv.fit(X_train, Y_train)  
  
[3]: print("Best hyperparameters : ", base_cv.best_params_)  
print("Accuracy : ", base_cv.best_score_)  
  
Best hyperparameters : {'n_neighbors': 10, 'algorithm': 'auto', 'p': 2}  
accuracy : 0.8442000000000001
```



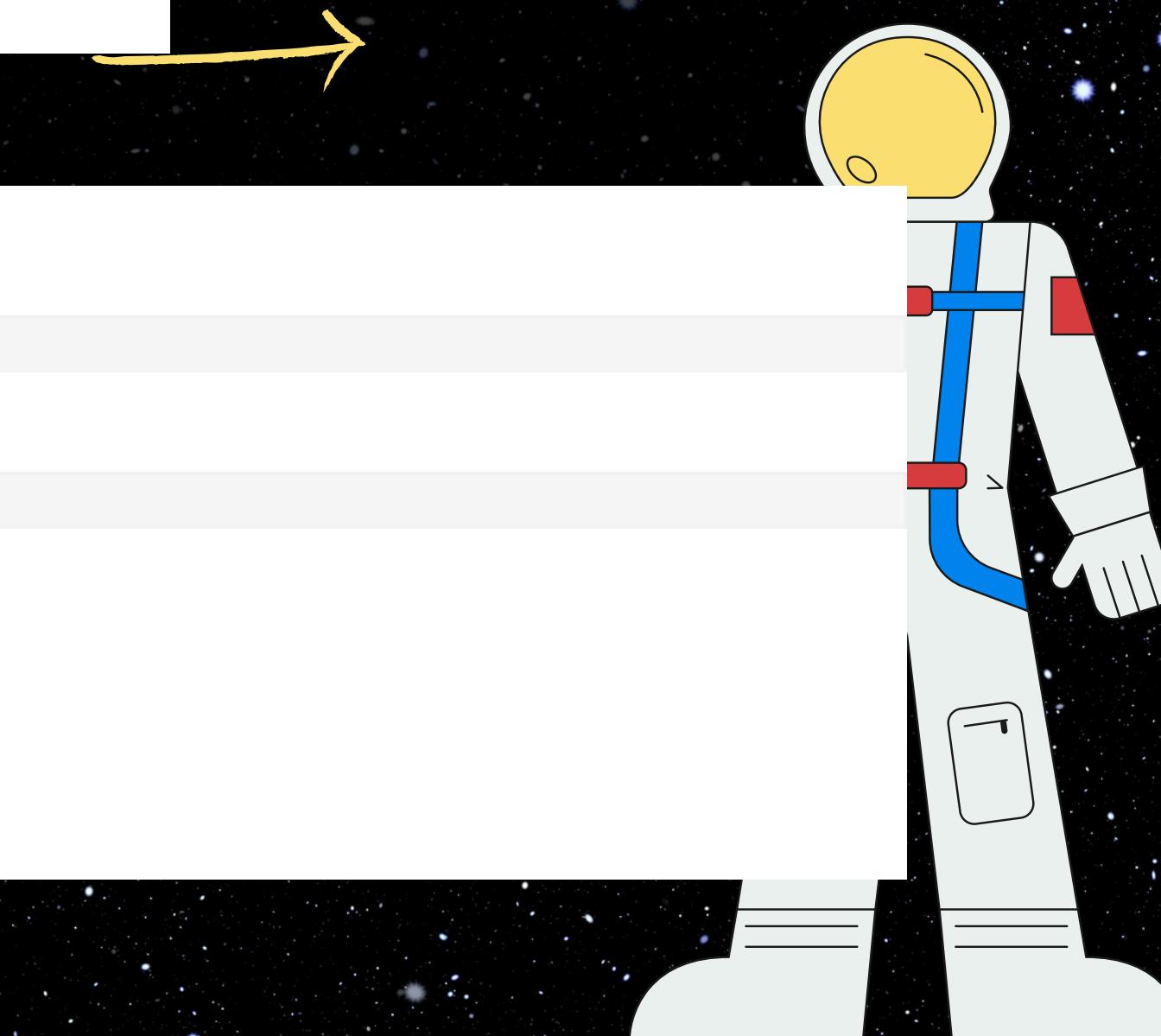
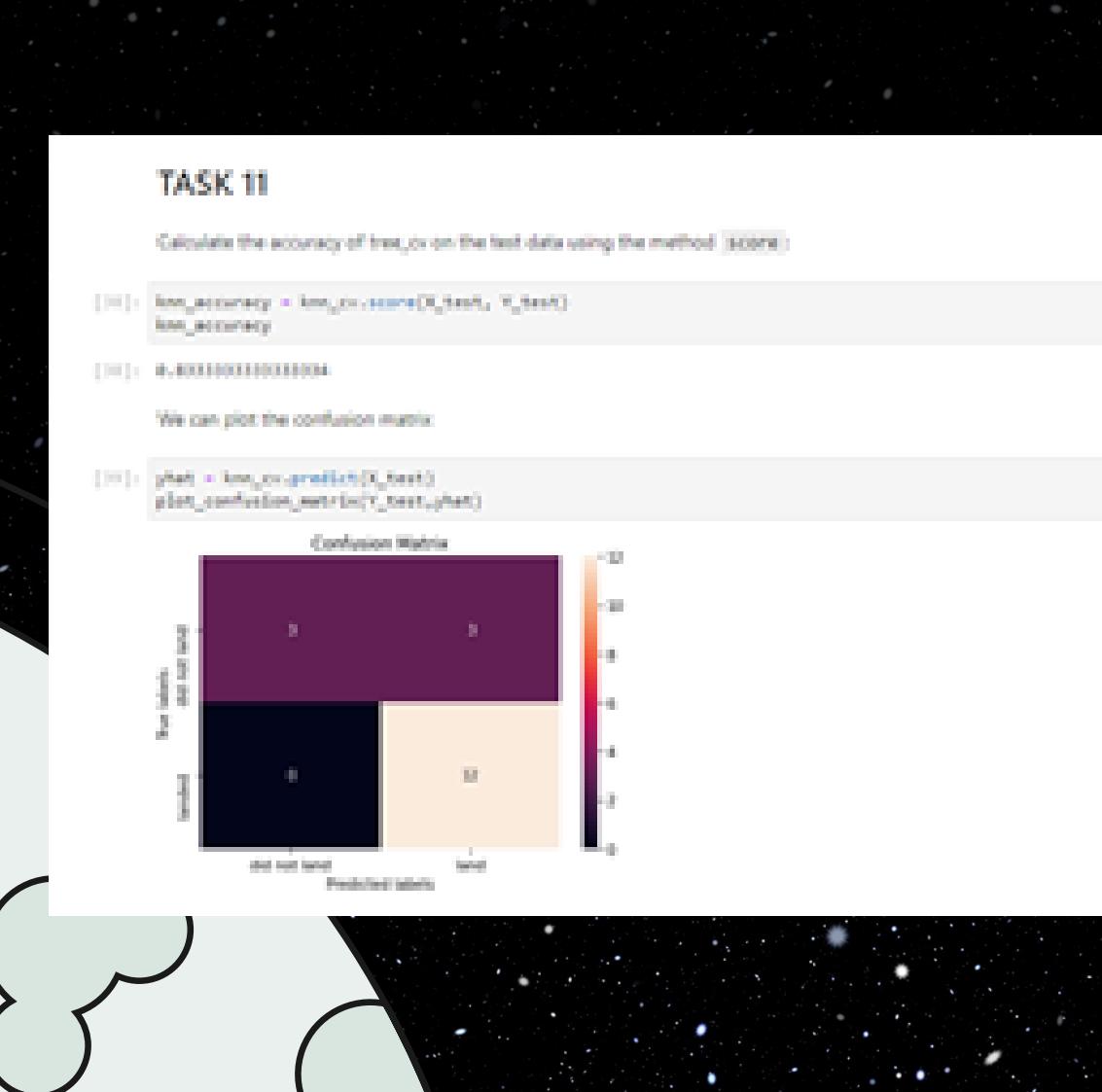
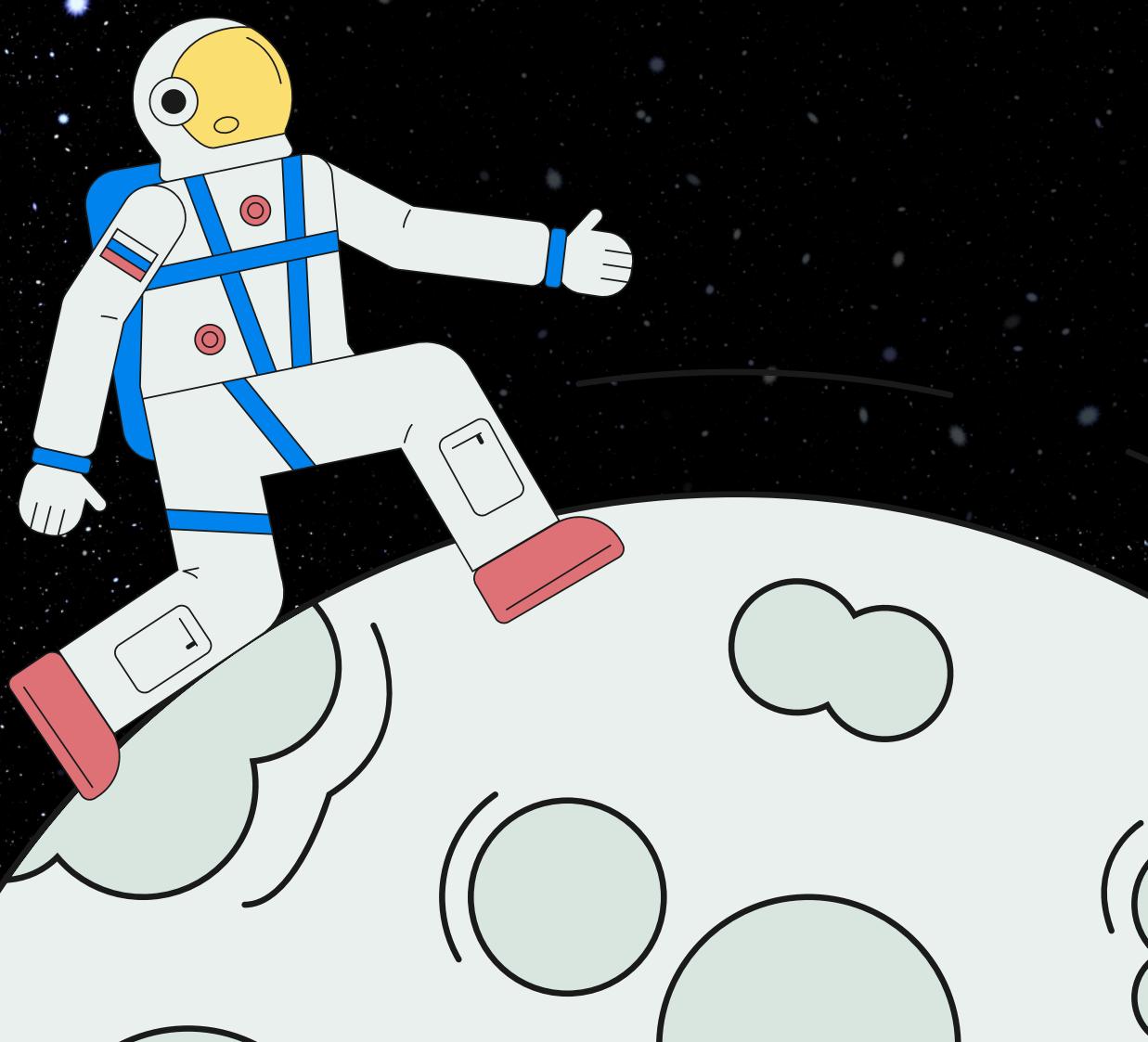
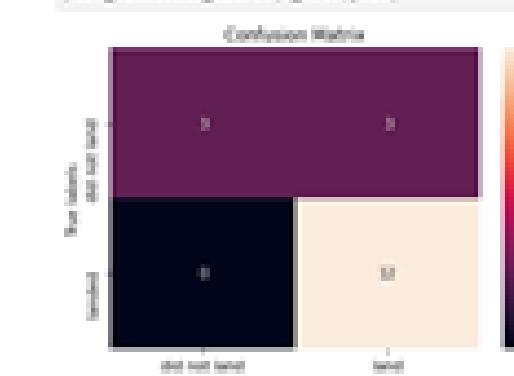
TASK 11

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

```
[1]: base_accuracy = base_cv.score(X_test, Y_test)  
base_accuracy
```

We can plot the confusion matrix:

```
[2]: pred = base_cv.predict(X_test)  
plot_confusion_matrix(base_cv, X_test, pred)
```



Results

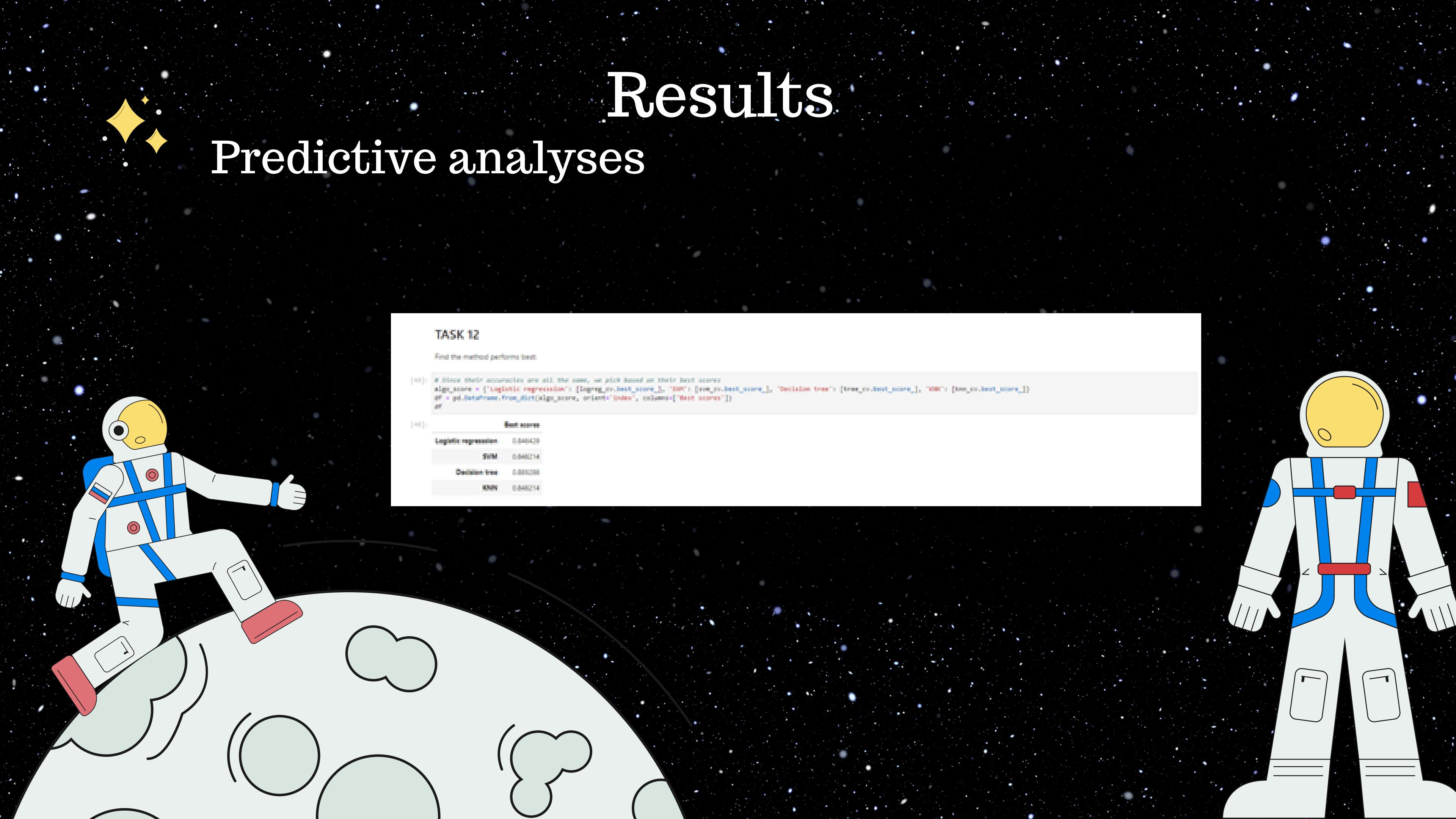
Predictive analyses

TASK 12

Find the method performs best.

```
[out]: As these stats scores are all the same, we pick based on their best scores
algos_scores = {"Logistic regression": [logreg_cv_cv_scores_1, "LR"], ["logreg_cv_cv_scores_2", "Logistic regression"], "Decision tree": [tree_cv_cv_scores_1, "DT"], ["tree_cv_cv_scores_2", "Decision tree"]]}
algos_scores["Decision tree", "best score"]
```

| | Best score |
|---------------------|------------|
| Logistic regression | 0.64627 |
| RF | 0.646214 |
| Decision tree | 0.646214 |
| KNN | 0.646214 |



Conclusion



Based on collected and wrangled data, we conclude that booster version F9 B5 B1048 - B1060 are the best booster versions for successful landing. Booster version F9 v1.1 B1012 and F9 v1.1 B1015 resulted in landing failures on drone ships. A greater success rate was observed ES-L1, SSO, HEO, and GEO. From the year 2013, significant improvements in success landing rates were observed. However, in 2017, a dip in success rates was observed.

