



IBM Developer
SKILLS NETWORK

Winning Space Race with Data Science

Jay Huang
08/10/2022



Outline

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

Executive Summary

- Summary of methodologies
- Summary of all results:
 - By using different classification models such as Logsitic Regression, Support Vector Machine, Decision Tree, and K nearest neighbor.
 - Gaining accuracy of 83.33%, 83.33%, 77.8%, 83.33% respectively.

Introduction

- Project background and context
- Problems you want to find answers

Section 1

Methodology

Methodology

Executive Summary

- Data collection methodology:
 - Data were collected by using wget with URL
- Perform data wrangling
 - Data were processed by removing null value, using standard scaler to normalize them, as well as encode categorical data
- Perform exploratory data analysis (EDA) using visualization and SQL
- Perform interactive visual analytics using Folium and Plotly Dash
- Perform predictive analysis using classification models
 - Built classification models by applying packages from Sklearn and tune model using GridSearchCV. Using AUC-ROC curve, confusion matrix, and F1 score to evaluates

Data Collection

- Describe how data sets were collected.
- You need to present your data collection process use key phrases and flowcharts

```
spacex_csv_file = wget.download('ht
```

```
data = pd.read_csv("https://cf-courses-dat
```

```
df=pd.read_csv("https://cf-courses-data.s3.us.cloud-object-storage.ap
```

```
1 %sql sqlite:///my_data1.db
UsageError: Line magic function

1 import pandas as pd
2 df = pd.read_csv("https://c
3 df.to_sql("SPACEXTBL", con,
```

Data Collection – SpaceX API

- Present your data collection with SpaceX REST calls using key phrases and flowcharts
- Add the GitHub URL of the completed SpaceX API calls notebook (must include completed code cell and outcome cell), as an external reference and peer-review purpose

FlightNumber	Date	BoosterVersion	PayloadMass	Orbit	LaunchSite	Outcome	Flights	GridFins	Reused	Legs	LandingPad	Block	Reus	
0	1	2010-06-04	Falcon 9	6104.959412	LEO	CCAFS SLC 40	None None	1	False	False	False	NaN	1.0	
1	2	2012-05-22	Falcon 9	525.000000	LEO	CCAFS SLC 40	None None	1	False	False	False	NaN	1.0	
2	3	2013-03-01	Falcon 9	677.000000	ISS	CCAFS SLC 40	None None	1	False	False	False	NaN	1.0	
3	4	2013-09-29	Falcon 9	500.000000	PO	VAFB SLC 4E	False Ocean	1	False	False	False	NaN	1.0	
4	5	2013-12-03	Falcon 9	3170.000000	GTO	CCAFS SLC 40	None None	1	False	False	False	NaN	1.0	
FlightNumber	PayloadMass	Flights	Block	ReusedCount	Orbit_ES-L1	Orbit_GEO	Orbit_GTO	Orbit_HEO	Orbit_ISS	...	Serial_B1058	S		
0	1.0	6104.959412	1.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0		
1	2.0	525.000000	1.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0		
2	3.0	677.000000	1.0	1.0	0.0	0.0	0.0	0.0	0.0	1.0	...	0.0		
3	4.0	500.000000	1.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0		
4	5.0	3170.000000	1.0	1.0	0.0	0.0	0.0	1.0	0.0	0.0	...	0.0		
...		
85	86.0	15400.000000	2.0	5.0	2.0	0.0	0.0	0.0	0.0	0.0	...	0.0		
86	87.0	15400.000000	3.0	5.0	2.0	0.0	0.0	0.0	0.0	0.0	...	1.0		
87	88.0	15400.000000	6.0	5.0	5.0	0.0	0.0	0.0	0.0	0.0	...	0.0		
88	89.0	15400.000000	3.0	5.0	2.0	0.0	0.0	0.0	0.0	0.0	...	0.0		
89	90.0	3681.000000	1.0	5.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0		

Data Collection - Scraping

- Present your web scraping process using key phrases and flowcharts
- Add the GitHub URL of the completed web scraping notebook, as an external reference and peer-review purpose

```
data = pd.read_csv("https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBM-DS0321EN-SkillsNetwork/datasets/dataset_part_2.csv")
```

```
spacex_csv_file = wget.download('https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBM-DS0321EN-SkillsNetwork/datasets/spacex_launch_geo.csv')
```

Data Wrangling

- Describe how data were processed
 - Built classification models by applying packages from Sklearn and tune model using GridSearchCV. Using AUC-ROC curve, confusion matrix, and F1 score to evaluate
- You need to present your data wrangling process using key phrases and flowcharts
- Add the GitHub URL of your completed data wrangling related notebooks, as an external reference and peer-review purpose

```
1 # Select relevant sub-columns: 'Launch Site', 'Lat(Latitude)', 'Long(Longitude)', 'class'
2 spacex_df = spacex_df[['Launch Site', 'Lat', 'Long', 'class']]
3 launch_sites_df = spacex_df.groupby(['Launch Site'], as_index=False).first()
4 launch_sites_df = launch_sites_df[['Launch Site', 'Lat', 'Long']]
5 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	-120.610746

```
1 # students get this
2 transform = preprocessing.StandardScaler()
3 X= transform.fit(X).transform(X)
4 X[0:5]
```

```
1 # HINT: Use get_dummies() function on the categorical columns
2 temp = pd.get_dummies(features[['Orbit', 'LaunchSite', 'LandingPad', 'Serial']])
```

```
1 features_one_hot = pd.concat([temp, features[['FlightNumber', 'PayloadMass', 'Flights', 'GridFins', 'Reus
2 features_one_hot.head()
```

	Orbit_ES-L1	Orbit_GEO	Orbit_GTO	Orbit_HEO	Orbit_ISS	Orbit_LEO	Orbit_MEO	Orbit_PO	Orbit_SO	Orbit_SSO	...	Serial
0	0	0	0	0	0	1	0	0	0	0	...	1
1	0	0	0	0	0	1	0	0	0	0	...	2
2	0	0	0	0	1	0	0	0	0	0	...	3
3	0	0	0	0	0	0	0	1	0	0	...	4
4	0	0	1	0	0	0	0	0	0	0	...	5

EDA with Data Visualization

- Summarize what charts were plotted and why you used those charts
 - Flight number V.S. Pay load mass(**Scatter plot**) to find out the outcome of the launch
 - Success rate of each orbit (**Bar chart**) to see which orbit perform better.
 - Year V.S. success rate(**line chart**) to find out how our tech perform when time goes by
- Add the GitHub URL of your completed EDA with data visualization notebook, as an external reference and peer-review purpose

EDA with SQL

- Using bullet point format, summarize the SQL queries you performed
 - Find unique launch sites in the space mission: `%sql Select DISTINCT LAUNCH_SITE FROM SPACEXTBL`
 - Show 5 records where launch sites begin with string 'CCA': `SELECT * FROM SPACEXTBL WHERE LAUNCH_SITE LIKE 'CCA%' LIMIT 5`
 - Total payload mass: `SELECT SUM(PAYLOAD_MASS_KG_) FROM SPACEXTBL WHERE CUSTOMER = 'NASA (CRS)'`
 - Find certain payload mass based on different booster version
 - Ground pad where first landing were successfully achieved
 - Name of booster when payload mass was between a certain range
 - Number of total outcomes
 - Outcomes between certain date range

Build an Interactive Map with Folium

- Summarize what map objects such as markers, circles, lines, etc. you created and added to a folium map
 - Using circle and add child to certain location to mark where it is; using marker to mark the size of the icon and color
 - Using folium.Map to show the whole map of USA; using circle to mark the radius and color as well as creating child.
 - Using mousePosition to get coordinate for a mouse over a point on the map
 - Using folium.Marker to mark the distance between to locations

Build a Dashboard with Plotly Dash

- Summarize what plots/graphs and interactions you have added to a dashboard
- Explain why you added those plots and interactions
- Add the GitHub URL of your completed Plotly Dash lab, as an external reference and peer-review purpose

Predictive Analysis (Classification)

- Summarize how you built, evaluated, improved, and found the best performing classification model
 - First by splitting the data into training and testing set
 - Using models from Sklearn packages to train each training set and then find out their accuracy
 - By using metrics like confusion matrix, F1 score and AUC-ROC curve we can directly see how the model perform on the data set
 - Besides decision tree, the other three models (logistic regression, SVM, KNN) have the same accuracy

Results

- Exploratory data analysis results
- Interactive analytics demo in screenshots
- Predictive analysis results

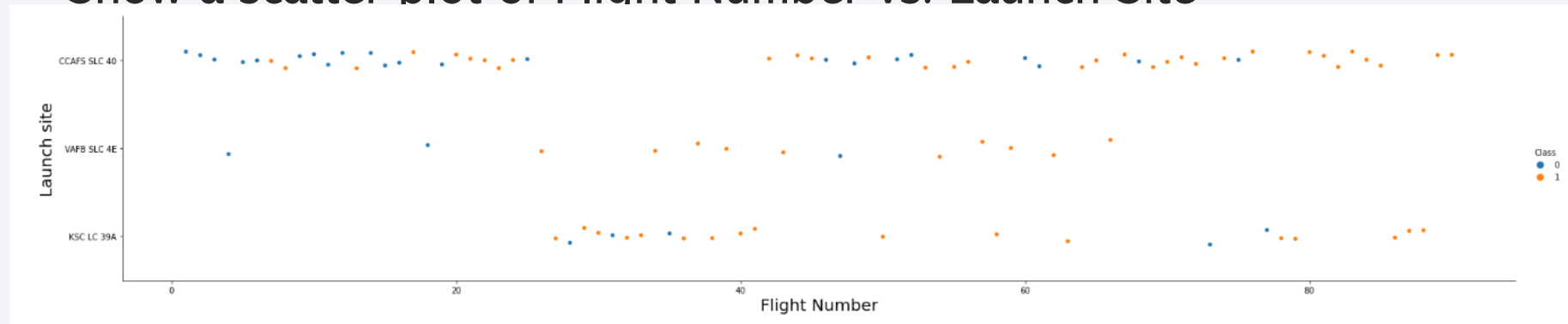
The background of the slide is an abstract composition. It features a dark blue base color. Overlaid on this are numerous diagonal streaks in shades of red and cyan. A faint, light blue grid pattern is also visible, particularly in the lower half of the image. The overall effect is dynamic and technological.

Section 2

Insights drawn from EDA

Flight Number vs. Launch Site

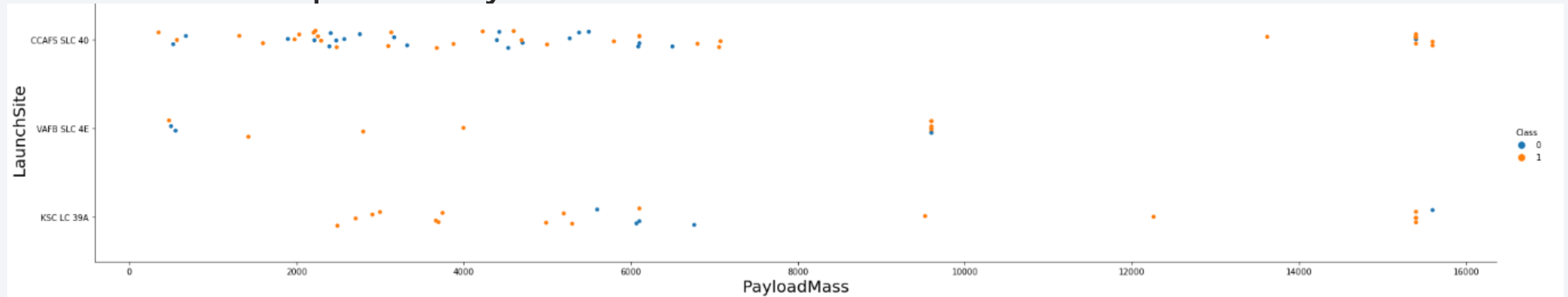
- Show a scatter plot of Flight Number vs. Launch Site



- The x-axis is flight number and y-axis is launch site. We can see that CCAFS SLC 40 site has the most flight number and VAFB SLC 4E has the least number of flights

Payload vs. Launch Site

- Show a scatter plot of Payload vs. Launch Site

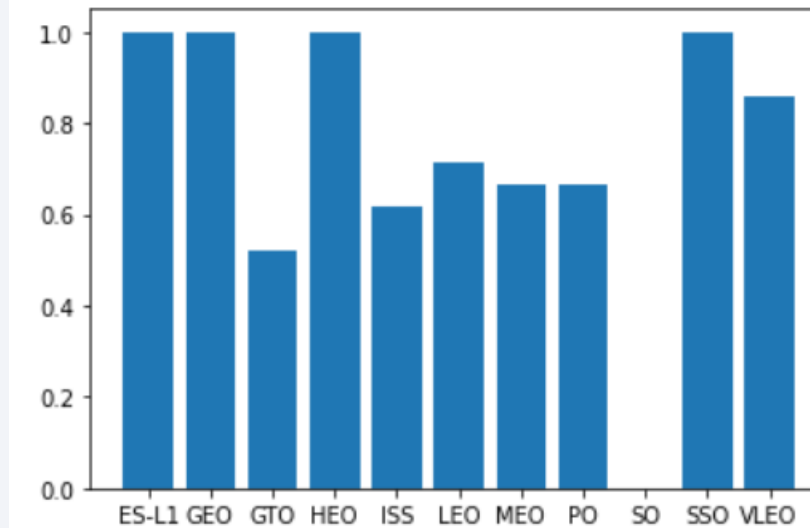


- Most of the payload are less than 8000, and still, CCAFS SLC 40 has most of the flights.

Success Rate vs. Orbit Type

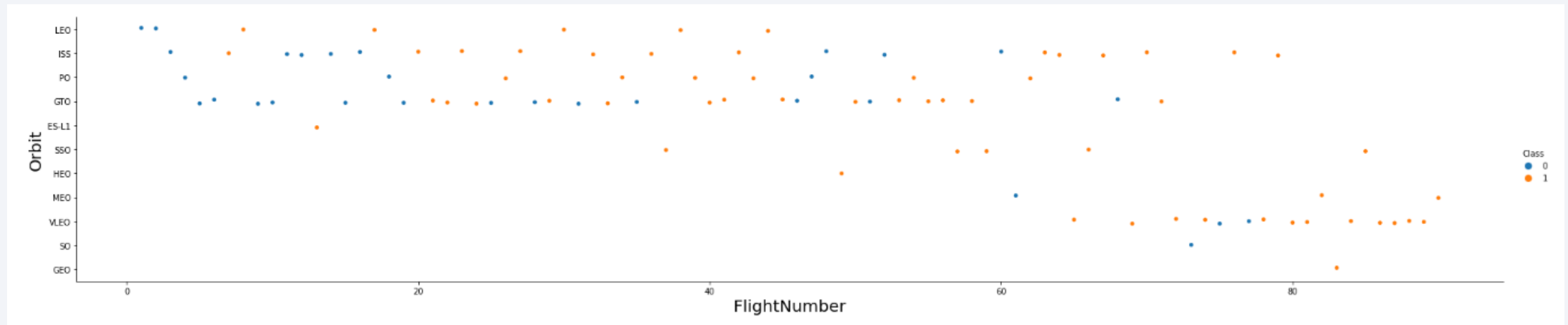
- Show a bar chart for the success rate of each orbit type

```
2 orbit = df[['Orbit', 'Class']].groupby('Orbit').mean()  
3  
4 plt.bar(orbit.index.values, orbit['Class'])  
5  
6 plt.show()
```



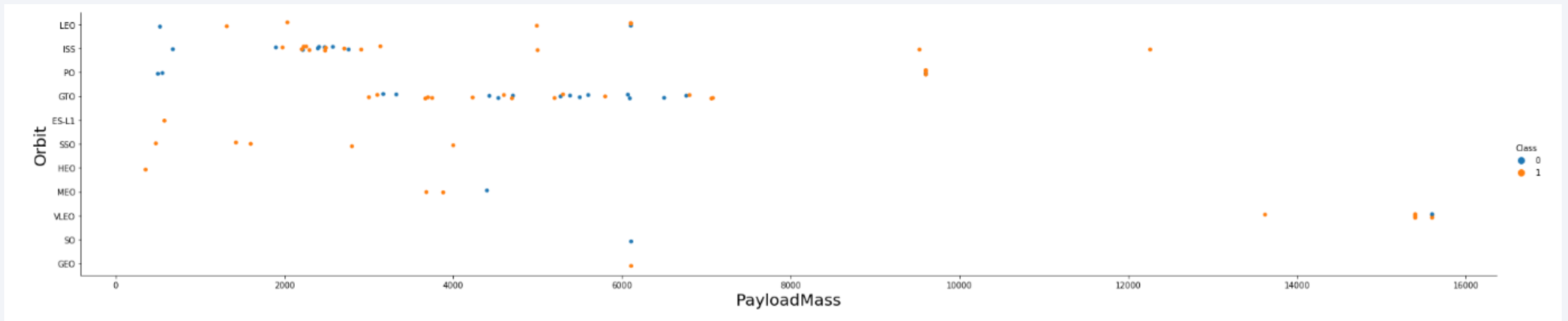
Flight Number vs. Orbit Type

- Show a scatter point of Flight number vs. Orbit type



Payload vs. Orbit Type

- Show a scatter point of payload vs. orbit type

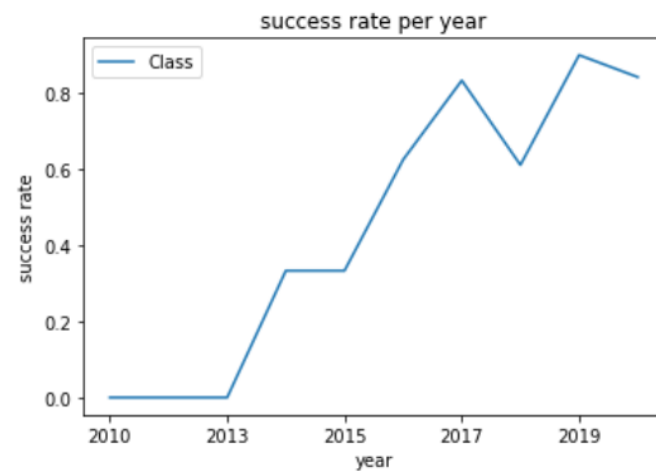


Launch Success Yearly Trend

- Show a line chart of yearly average success rate

```
2 df[['Class', 'year']].groupby('year').mean().plot()
3 plt.title('success rate per year')
4 plt.ylabel('success rate')
5 plt.xlabel('year')
```

Text(0.5, 0, 'year')



All Launch Site Names

- Find the names of the unique launch sites

```
%sql SELECT Distinct LAUNCH_SITE FROM SPACEXTBL
```

```
* ibm_db_sa://kcq64325:***@dashdb-txn-sbox-yp-dal  
Done.
```

launch_site

CCAFS LC-40

CCAFS SLC-40

KSC LC-39A

VAFB SLC-4E

Launch Site Names Begin with 'CCA'

- Find 5 records where launch sites begin with 'CCA'

• `%sql SELECT * FROM SPACESTBL WHERE LAUNCH_SITE LIKE 'CCA%' LIMIT 5`

* ibm_db_sa://kcq64325:***@dashdb-txn-sbox-yp-dal109-04.services.dal.ibm.com:50000/BLUDB
Done.

DATE	time_utc	booster_version	launch_site	payload	payload_mass_kg	orbit	customer	mission_outcome	landing_outcome
2010-06-04	18:45:00	F9 v1.0 B0003	CCAFS LC-40	None	0	LEO	SpaceX	Success	Failure (parachute)
2010-12-08	15:43:00	F9 v1.0 B0004	CCAFS LC-40	None	0	LEO (ISS)	NASA (COTS) NRO	Success	Failure (parachute)
2012-05-22	07:44:00	F9 v1.0 B0005	CCAFS LC-40	None	525	LEO (ISS)	NASA (COTS)	Success	No attempt
2012-10-08	00:35:00	F9 v1.0 B0006	CCAFS LC-40	None	500	LEO (ISS)	NASA (CRS)	Success	No attempt
2013-03-01	15:10:00	F9 v1.0 B0007	CCAFS LC-40	None	677	LEO (ISS)	NASA (CRS)	Success	No attempt

Total Payload Mass

- Calculate the total payload carried by boosters from NASA

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

```
* ibm_db_sa://kcq64325:***@dashdb-txn-sbox-yp-dal109-04.services.dal.ibmcloud.com:54321:
Done.
```

```
1
```

```
45596
```

Average Payload Mass by F9 v1.1

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

```
%sql SELECT AVG(PAYLOAD_MASS_KG_) FROM SPACEXTBL WHERE BOOSTER_VERSION='F9 v1.1'
```

```
* ibm_db_sa://kcq64325:***@dashdb-txn-sbox-yp-dal109-04.services.dal.ibmcloud.com:50000
Done.
```

```
1
```

```
2928.400000
```

First Successful Ground Landing Date

- Find the dates of the first successful landing outcome on ground pad

```
%sql SELECT min(DATE) FROM SPACEXTBL WHERE LANDING__OUTCOME='Success (ground pad)'  
  
* ibm_db_sa://kcq64325:***@dashdb-txn-sbox-yp-dal09-04.services.dal.ibmcloud.com:500  
Done.  
  
1  
2015-12-22
```


Successful Drone Ship Landing with Payload between 4000 and 6000

- List the names of boosters which have successfully landed on drone ship and had payload mass greater than 4000 but less than 6000

```
%sql SELECT BOOSTER_VERSION FROM SPACEXTBL WHERE PAYLOAD_MASS_KG_ between 4000 and 6000 AND LANDING__OUTCOME='Success (drone ship)'
```

```
* ibm_db_sa://kcq64325:***@dashdb-txn-sbox-yp-dal109-04.services.dal.ibmcloud.net:50000/BLUDB  
Done.
```

booster_version

F9 FT B1022

F9 FT B1026

F9 FT B1021.2

F9 FT B1031.2

Total Number of Successful and Failure Mission Outcomes

- Calculate the total number of successful and failure mission outcomes

```
%sql SELECT COUNT(*) FROM SPACEXTBL WHERE MISSION_OUTCOME LIKE '%Success%' OR MISSION_OUTCOME LIKE '%Failure%'
```

```
* ibm_db_sa://kcq64325:***@dashdb-txn-sbox-yp-dal09-04.services.dal.ibmcloud.net:50000/BLUDB  
Done.
```

```
1
```

```
101
```

Boosters Carried Maximum Payload

- List the names of the booster which have carried the maximum payload mass

```
%sql SELECT BOOSTER_VERSION FROM SPACEXTBL WHERE PAYLOAD_MASS_KG_ = (SELECT MAX(PAYLOAD_MASS_KG_) FROM SPACEXTBL)
```

```
* ibm_db_sa://kcq64325:***@dashdb-txn-sbox-yp-da109-04.services.dal1.bluemix.net:50000/BLUDB  
Done.
```

booster_version

F9 B5 B1048.4

F9 B5 B1049.4

F9 B5 B1051.3

F9 B5 B1056.4

F9 B5 B1048.5

F9 B5 B1051.4

F9 B5 B1049.5

F9 B5 B1060.2

F9 B5 B1058.3

F9 B5 B1051.6

F9 B5 B1060.3

F9 B5 B1049.7

2015 Launch Records

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

```
%sql SELECT TO_CHAR(TO_DATE(MONTH("DATE"), 'MM'), 'MONTH') AS MONTH_NAME, \
LANDING__OUTCOME AS LANDING__OUTCOME, \
BOOSTER_VERSION AS BOOSTER_VERSION, \
LAUNCH_SITE AS LAUNCH_SITE \
FROM SPACEXTBL WHERE LANDING__OUTCOME = 'Failure (drone ship)' AND "DATE" LIKE '%2015%'
```

```
* ibm_db_sa://kcq64325:***@dashdb-txn-sbox-yp-dal09-04.services.dal.ibm.com:50000/BLUDE
Done.
```

month_name	landing_outcome	booster_version	launch_site
JANUARY	Failure (drone ship)	F9 v1.1 B1012	CCAFS LC-40
APRIL	Failure (drone ship)	F9 v1.1 B1015	CCAFS LC-40

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)) between the date 2010-06-04 and 2017-03-20, in

```
%sql SELECT "DATE", COUNT(LANDING__OUTCOME) as COUNT FROM SPACEXTBL \
WHERE "DATE" BETWEEN '2010-06-04' and '2017-03-20' AND LANDING__OUTCOME LIKE '%Success%' \
GROUP BY "DATE" \
ORDER BY COUNT(LANDING__OUTCOME) DESC
```

```
* ibm_db_sa://kcq64325:***@dashdb-txn-sbox-yp-da109-04.services.dal1.bluemix.net:50000/BLUDB
Done.
```

DATE	COUNT
2015-12-22	1
2016-04-08	1
2016-05-06	1
2016-05-27	1
2016-07-18	1
2016-08-14	1
2017-01-14	1
2017-02-19	1

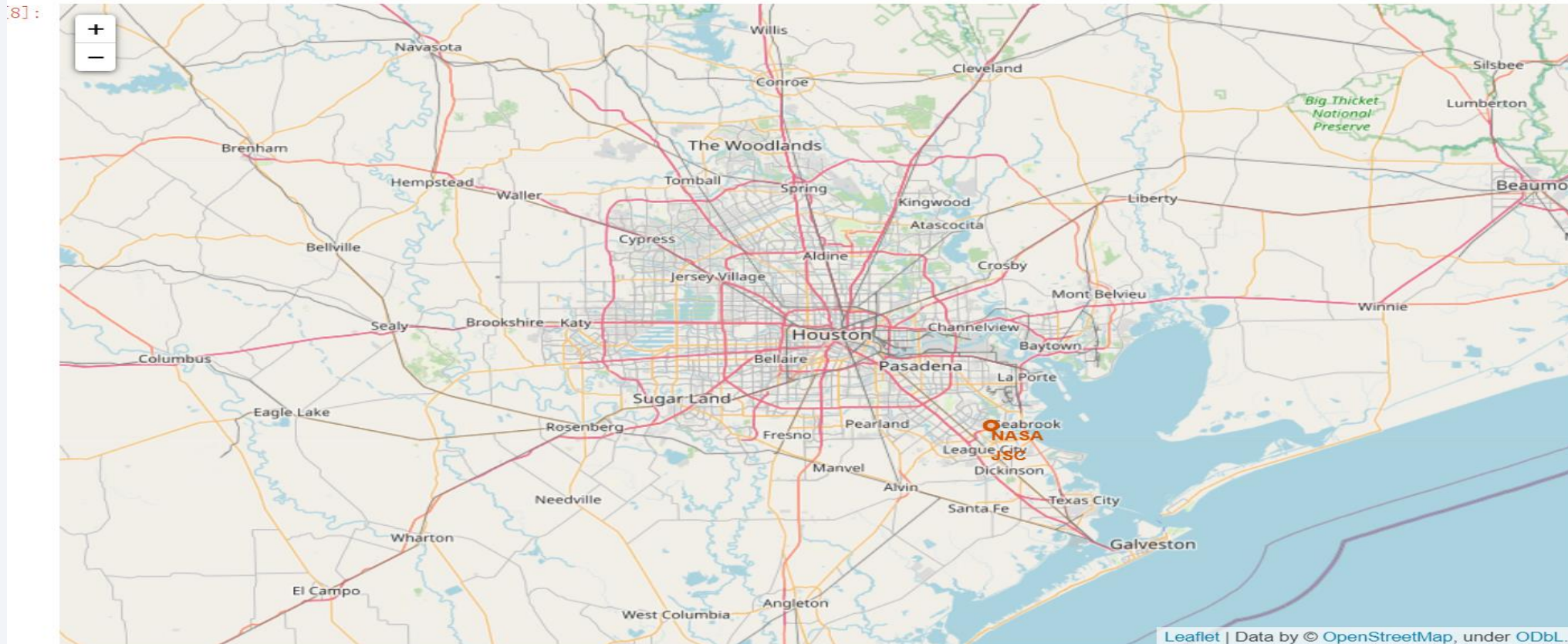
A satellite view of Earth from space, showing the curvature of the planet and city lights at night. The background is a deep blue gradient.

Section 3

Launch Sites Proximities Analysis

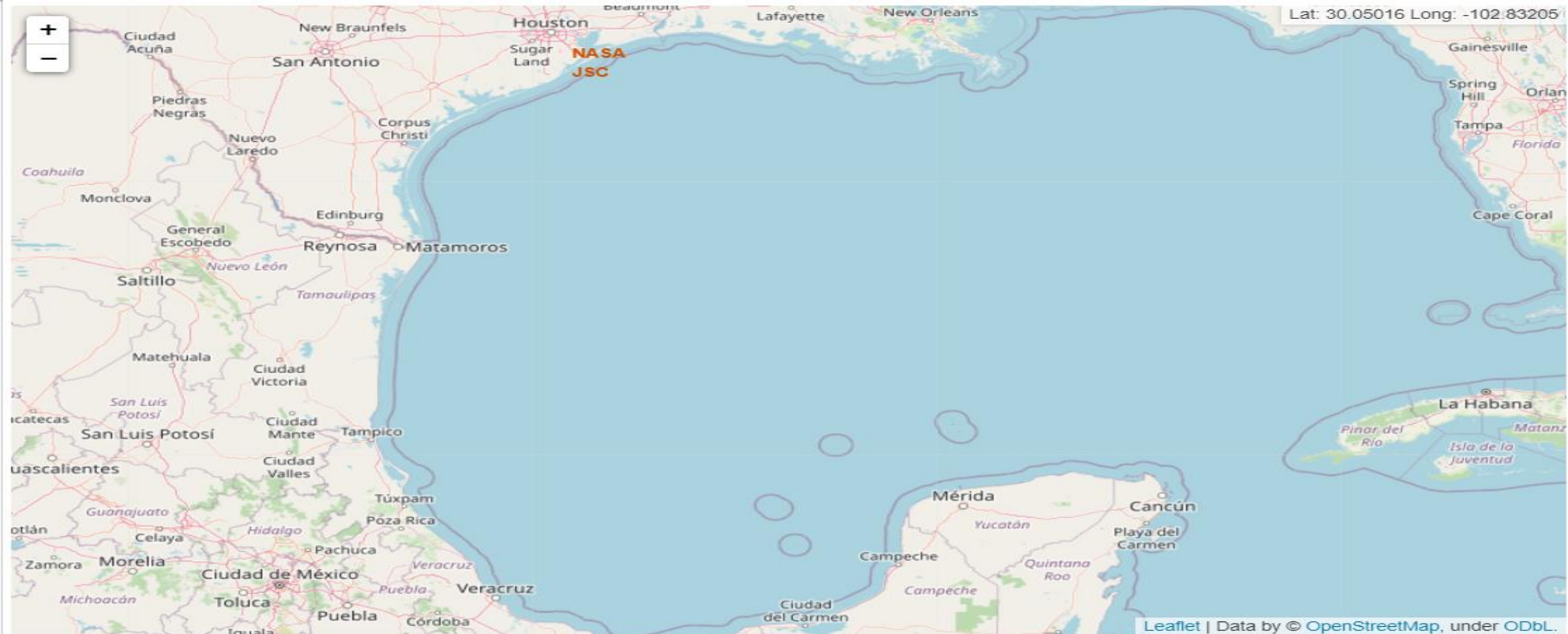
Houston NASA map

```
2 circle = folium.Circle(nasa_coordinate, radius=1000, color='#d35400', fill=True).add_child(folium.Popup('NASA Johnson Space Center'))
3 # Create a blue circle at NASA Johnson Space Center's coordinate with a icon showing its name
4 marker = folium.map.Marker(
5     nasa_coordinate,
6     # Create an icon as a text label
7     icon=DivIcon(
8         icon_size=(20, 20),
9         icon_anchor=(0, 0),
10        html='<div style="font-size: 12; color:#d35400;"><b>%s</b></div>' % 'NASA JSC',
11    )
12 )
13 site_map.add_child(circle)
14 site_map.add_child(marker)
```



Distance Calculation Map

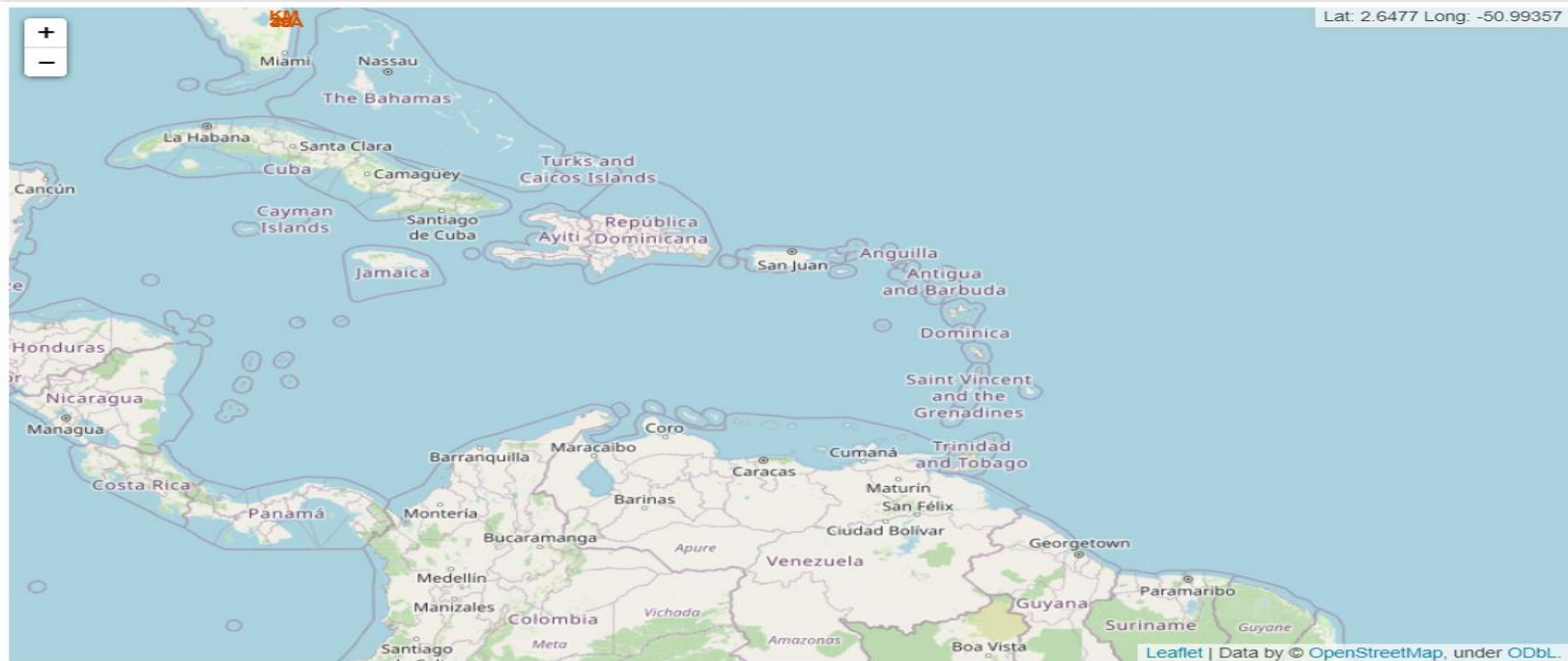
```
2 formatter = "function(num) {return L.Util.formatNum(num, 5)};"
3 mouse_position = MousePosition(
4     position=' topright',
5     separator=' Long: ',
6     empty_string=' NaN',
7     lng_first=False,
8     num_digits=20,
9     prefix=' Lat:',
10    lat_formatter=formatter,
11    lng_formatter=formatter,
12 )
13
14 site_map.add_child(mouse_position)
15 site_map
```



Southern map

```
23 )
24 site_map.add_child(distance_marker)
25 # closest railroad line
26 coordinates = [[launch_site_lat, launch_site_lon], closest_railroad]
27 lines=folium.PolyLine(locations=coordinates, weight=1)
28 site_map.add_child(lines)
29
30 # closest city marker
31 distance_marker = folium.Marker(
32     closest_city,
33     icon=DivIcon(
34         icon_size=(20,20),
35         icon_anchor=(0,0),
36         html='<div style="font-size: 12; color:#d35400;"><b>%s</b></div>' % "{:10.2f} KM".format(distance_city),
37     )
38 )
39 site_map.add_child(distance_marker)
40 # closest city line
41 coordinates = [[launch_site_lat, launch_site_lon], closest_city]
42 lines=folium.PolyLine(locations=coordinates, weight=1)
43 site_map.add_child(lines)
```

Out[23]:

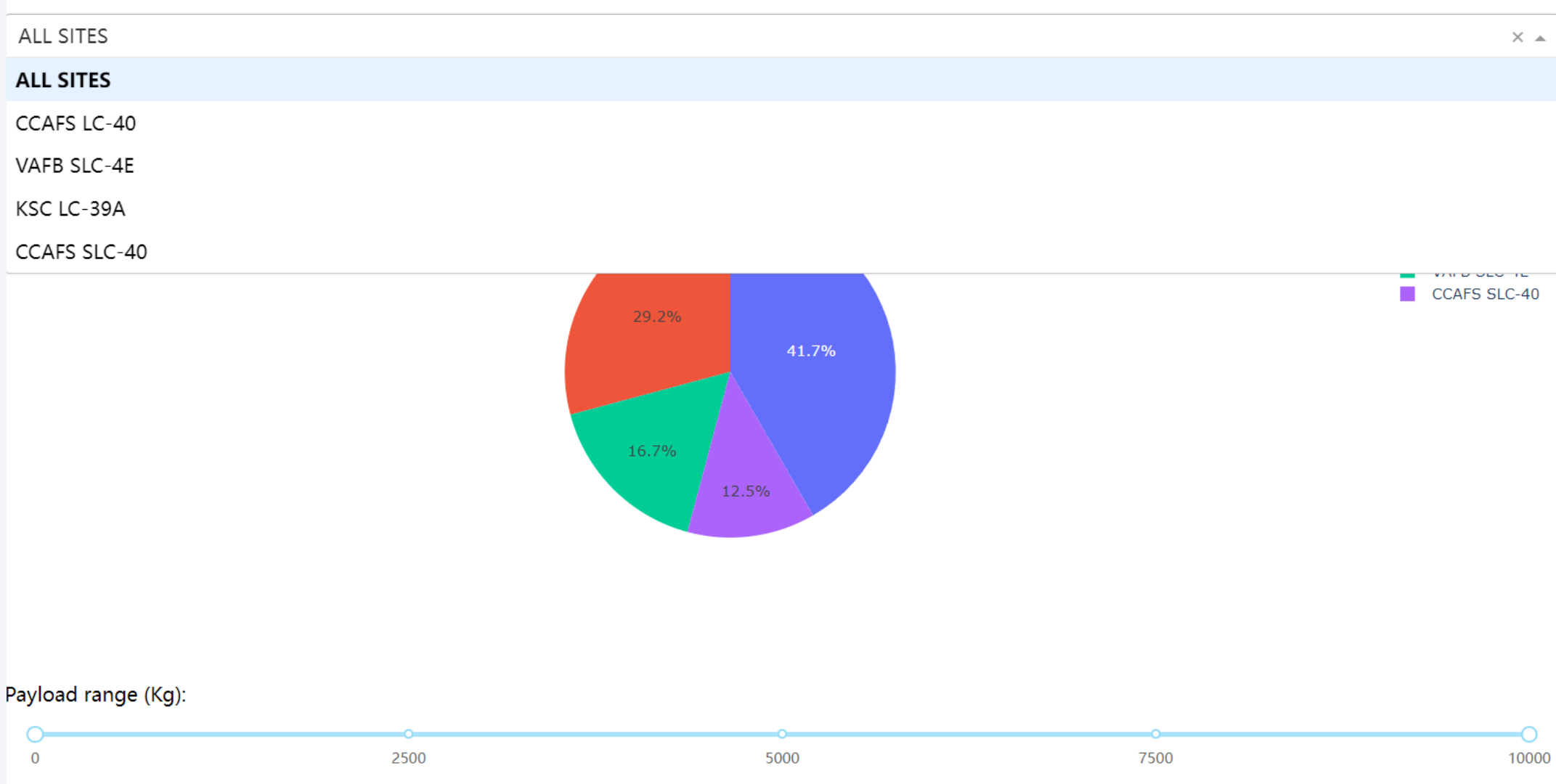


The background of the slide is a close-up, artistic photograph of a printed circuit board (PCB). The board is dark, and the intricate circuit traces are highlighted in a vibrant, glowing red. Numerous small, circular components, likely solder joints or micro-components, are visible along the traces, some of which also appear to be glowing. The overall effect is a high-tech, digital aesthetic.

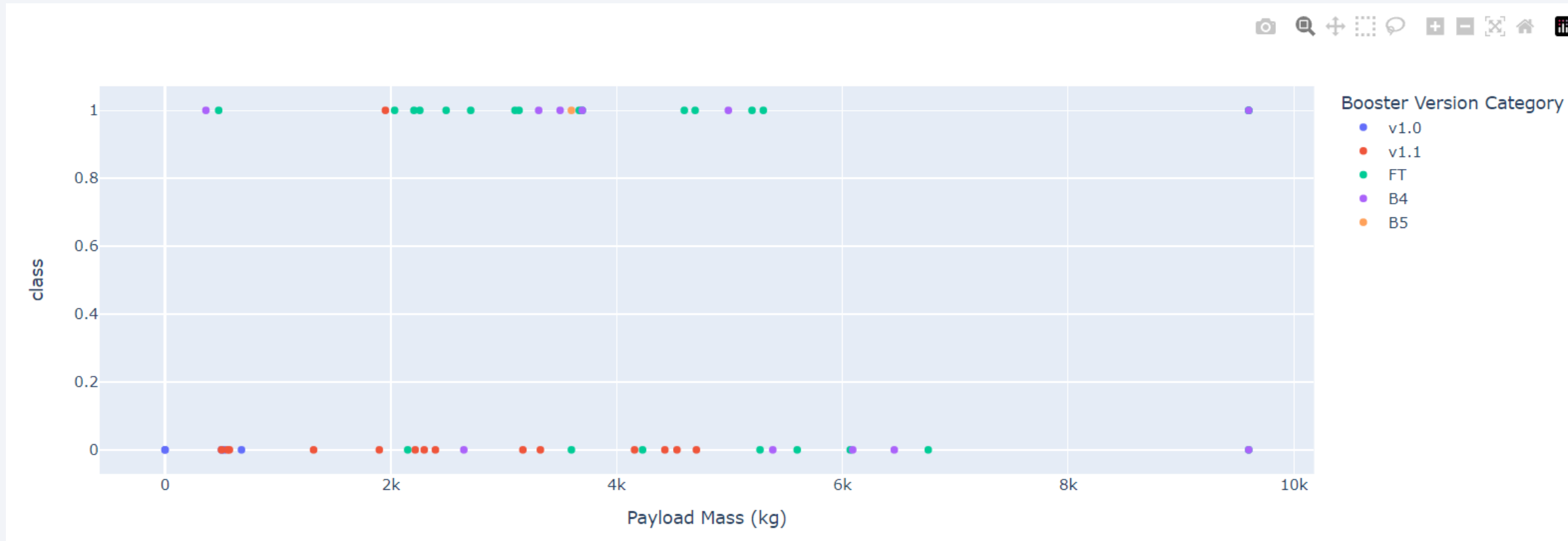
Section 4

Build a Dashboard with Plotly Dash

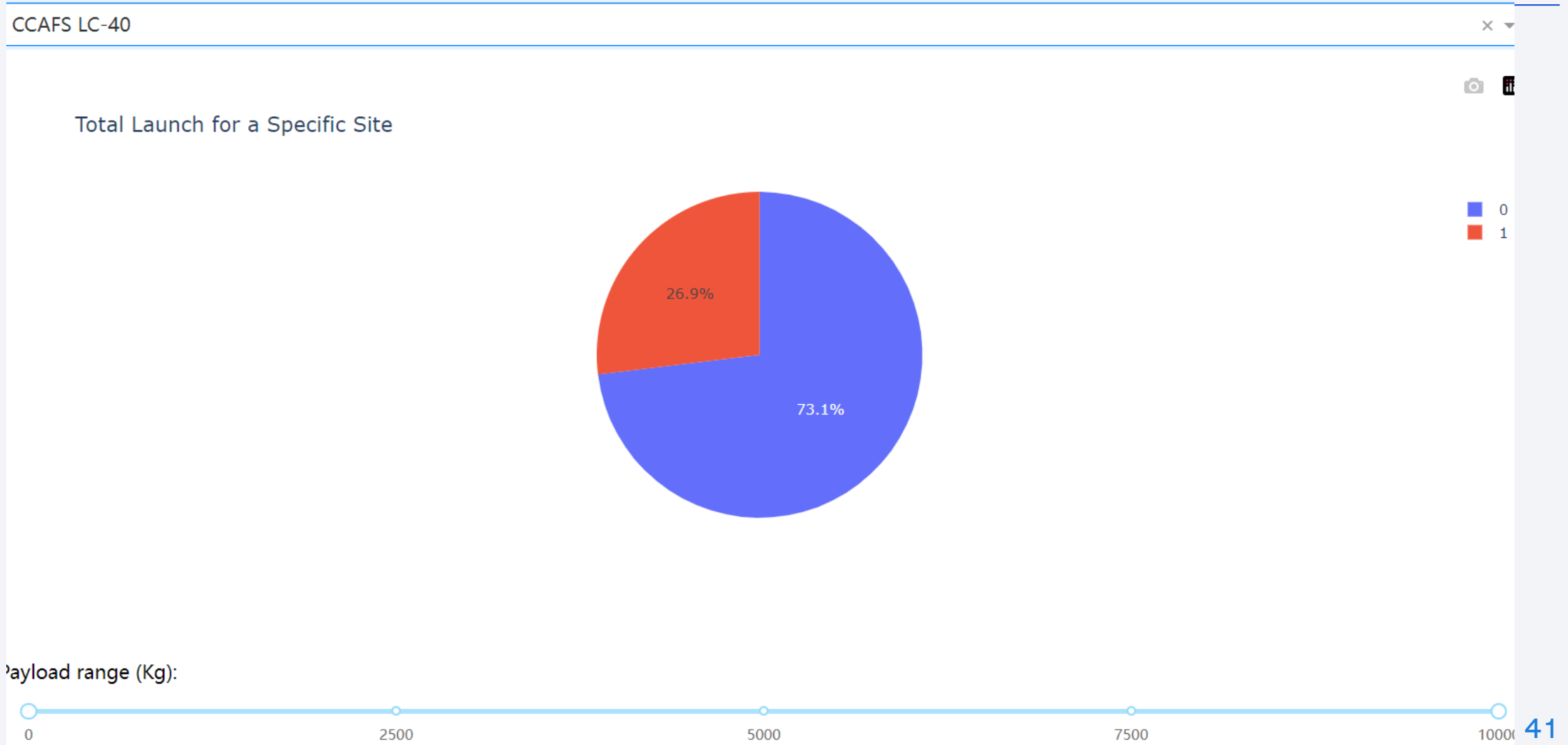
Sites V.S. payload range



Payload V.S. Class



Site that has the most flights



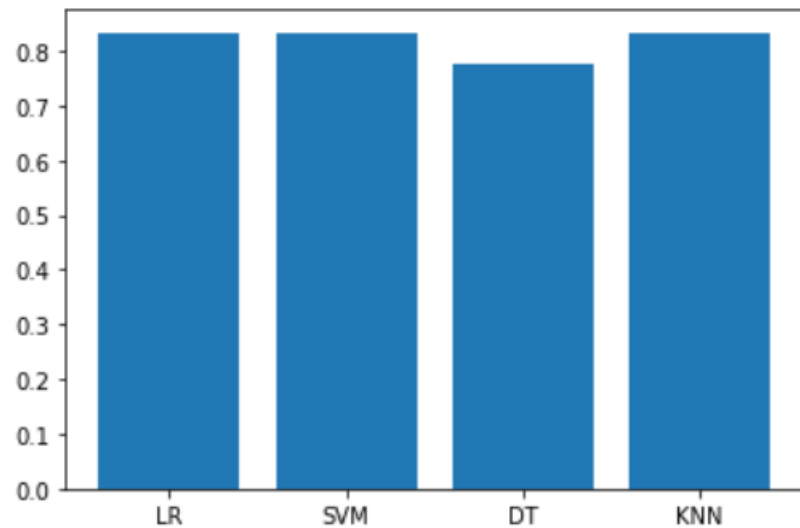
Section 5

Predictive Analysis (Classification)

Classification Accuracy

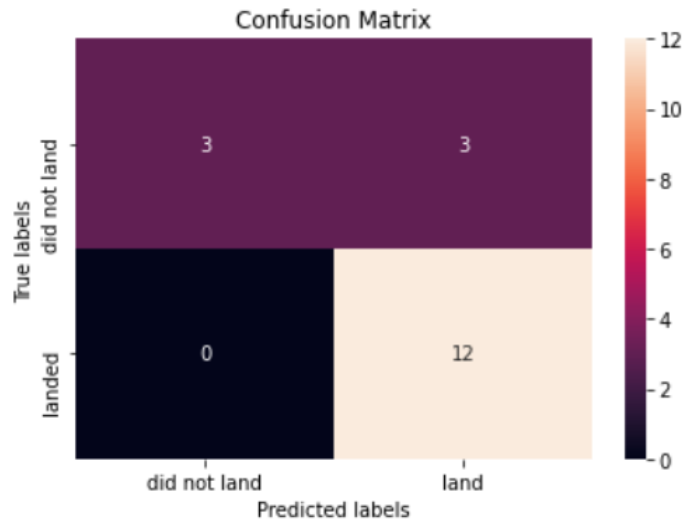
```
1 accuracy_lst = [logreg_cv.score(X_test, Y_test), svm_cv.score(X_test, Y_test), tree_cv.score(X_test, Y_test), knn_cv.score(X_test, Y_test)]  
2 name_lst = ['LR', 'SVM', 'DT', 'KNN']  
3 plt.bar(name_lst, accuracy_lst)
```

<BarContainer object of 4 artists>

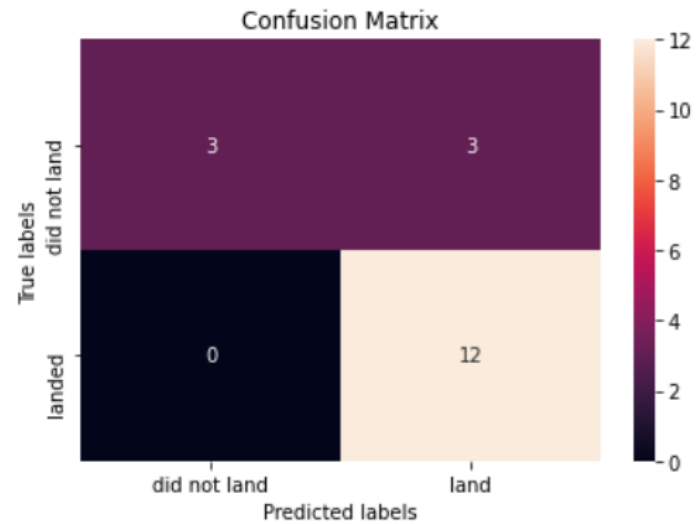


Confusion Matrix

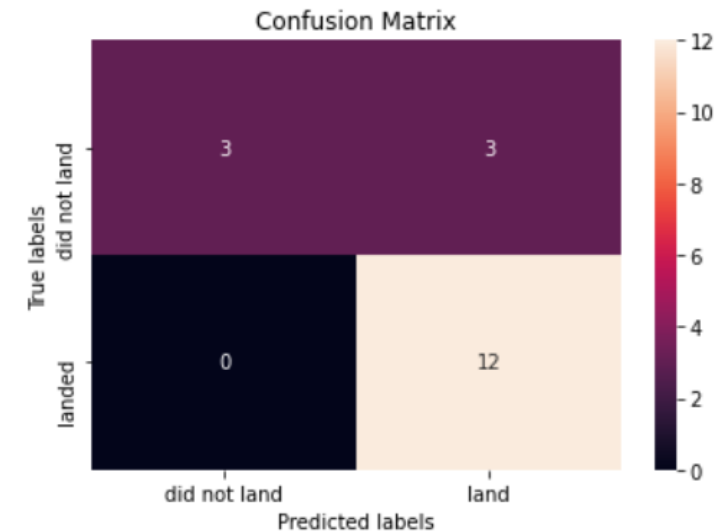
```
1 #KNN  
2 yhat = knn_cv.predict(X_test)  
3 plot_confusion_matrix(Y_test, yhat)
```



```
1 #SVM  
2 yhat=svm_cv.predict(X_test)  
3 plot_confusion_matrix(Y_test, yhat)
```



```
1 #Logistic Regression  
2 yhat=logreg_cv.predict(X_test)  
3 plot_confusion_matrix(Y_test, yhat)
```



Conclusions

- Point 1
- Point 2
- Point 3
- Point 4
- ...

Appendix

- Include any relevant assets like Python code snippets, SQL queries, charts, Notebook outputs, or data sets that you may have created during this project

Thank you!

