



IBM Developer
SKILLS NETWORK

Winning Space Race with Data Science

Seyed Mohammad Ahmadi
May 26 2023



Outline

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

Executive Summary

- Introduction
- Summary of methodologies
- Data Collection & Data Wrangling
- EDA & Interactive Visual Analytics Methodology
- Predictive Analysis Methodology
- EDA with Visualization Results
- EDA With SQL Results
- Interactive Map with Folium Results
- Plotly Dashboard Results
- Predictive Analysis (Classification) Results
- Summary of all results

Introduction

- **Project background and context**

SpaceX advertises Falcon 9 rocket launches on its website with a cost of 62 million dollars; other providers cost upwards of 165 million dollars each, much of the savings is because SpaceX can reuse the first stage.

- **Problems you want to find answers**

If we can determine if the first stage will land, we can determine the cost of a launch

Section 1

Methodology

Methodology

Executive Summary

- Data collection methodology:
 - SpaceX Rest API
 - Web Scrapping from SpaceX Wikipedia webpage
- Perform data wrangling
 - Hot Encoding data fields
 - Data Cleaning: Correcting and Cleaning null values & irrelevant columns
- Perform exploratory data analysis (EDA) using visualization and SQL
- Perform interactive visual analytics using Folium and Plotly Dash
- Perform predictive analysis using classification models
 - Different algorithm such as: Logistic Regression, K-Nearest Neighbors, Support Vector Machine, Decision Tree model have been deployed and evaluated to find the best method.

Data Collection

There are two methods that data are collected:

1

Used SpaceX
Rest API

Returned
SpaceX data

Normalized data
int a csv file

Data
Consolidation &
Wrangling

2

Got HTML
Response from
Wikipedia webpage

Used BeautifulSoup
to extract Data

Normalized data int
a csv file

Data Consolidation
& Wrangling

Data Collection – SpaceX API

- Data collection with SpaceX REST calls

Out[35]:

	FlightNumber	Date	BoosterVersion	PayloadMass	Orbit	LaunchSite	Outcome	Flights	GridFins	Reused	Legs	La
4	1	2010-06-04	Falcon 9	NaN	LEO	CCSFS SLC 40	None None	1	False	False	False	
5	2	2012-05-22	Falcon 9	525.0	LEO	CCSFS SLC 40	None None	1	False	False	False	
6	3	2013-03-01	Falcon 9	677.0	ISS	CCSFS SLC 40	None None	1	False	False	False	
7	4	2013-09-29	Falcon 9	500.0	PO	VAFB SLC 4E	False Ocean	1	False	False	False	
8	5	2013-12-03	Falcon 9	3170.0	GTO	CCSFS SLC 40	None None	1	False	False	False	
...	
89	86	2020-09-03	Falcon 9	15600.0	VLEO	KSC LC 39A	True ASDS	2	True	True	True	5e9e3032383ecb6b

- <https://github.com/Mamad66Ahmadi/IBM-Data-Science-Capstone/blob/Module1/jupyter-labs-spacex-data-collection-api.ipynb>

Data Collection - Scraping

- Web scraping process

2020 [edit]

In late 2019, Gwynne Shotwell stated that SpaceX hoped for as many as 24 launches for Starlink satellites in 2020,^[490] in addition to 14 or 15 non-Starlink launches. At 26 launches, 13 of which for Starlink satellites, Falcon 9 had its most prolific year, and Falcon rockets were second most prolific rocket family of 2020, only behind China's Long March rocket family.^[491]

[hide] Flight No.	Date and time (UTC)	Version, Booster ^[b]	Launch site	Payload ^[c]	Payload mass	Orbit	Customer	Launch outcome	Booster landing
78	7 January 2020, 02:19:21 ^[492]	F9 B5 △ B1049.4	CCAFS, SLC-40	Starlink 2 v1.0 (60 satellites)	15,600 kg (34,400 lb) ^[5]	LEO	SpaceX	Success	Success (drone ship)
Third large batch and second operational flight of Starlink constellation. One of the 60 satellites included a test coating to make the satellite less reflective, and thus less likely to interfere with ground-based astronomical observations. ^[493]									
79	19 January 2020, 15:30 ^[494]	F9 B5 △ B1046.4	KSC, LC-39A	Crew Dragon in-flight abort test ^[495] (Dragon C205.1)	12,050 kg (26,570 lb)	Sub-orbital ^[496]	NASA (CTS) ^[497]	Success	No attempt
An atmospheric test of the Dragon 2 abort system after Max Q. The capsule fired its SuperDraco engines, reached an apogee of 40 km (25 mi), deployed parachutes after reentry, and splashed down in the ocean 31 km (19 mi) downrange from the launch site. The test was previously slated to be accomplished with the Crew Dragon Demo-1 capsule, ^[498] but that test article exploded during a ground test of SuperDraco engines on 20 April 2019. ^[419] The abort test used the capsule originally intended for the first crewed flight. ^[499] As expected, the test was successful. The capsule was recovered by the USS <i>Thetis</i> (MCMC-103) and returned to the launch site. The second stage had a mass simulator in place of its engine.									
80	29 January 2020, 14:07 ^[501]						SpaceX	Success	Success (drone ship)
Third operational and flight of Starlink constellation. The capsule was recovered by the USS <i>Thetis</i> (MCMC-103) and returned to the launch site.									
17 February 2020,									Failure

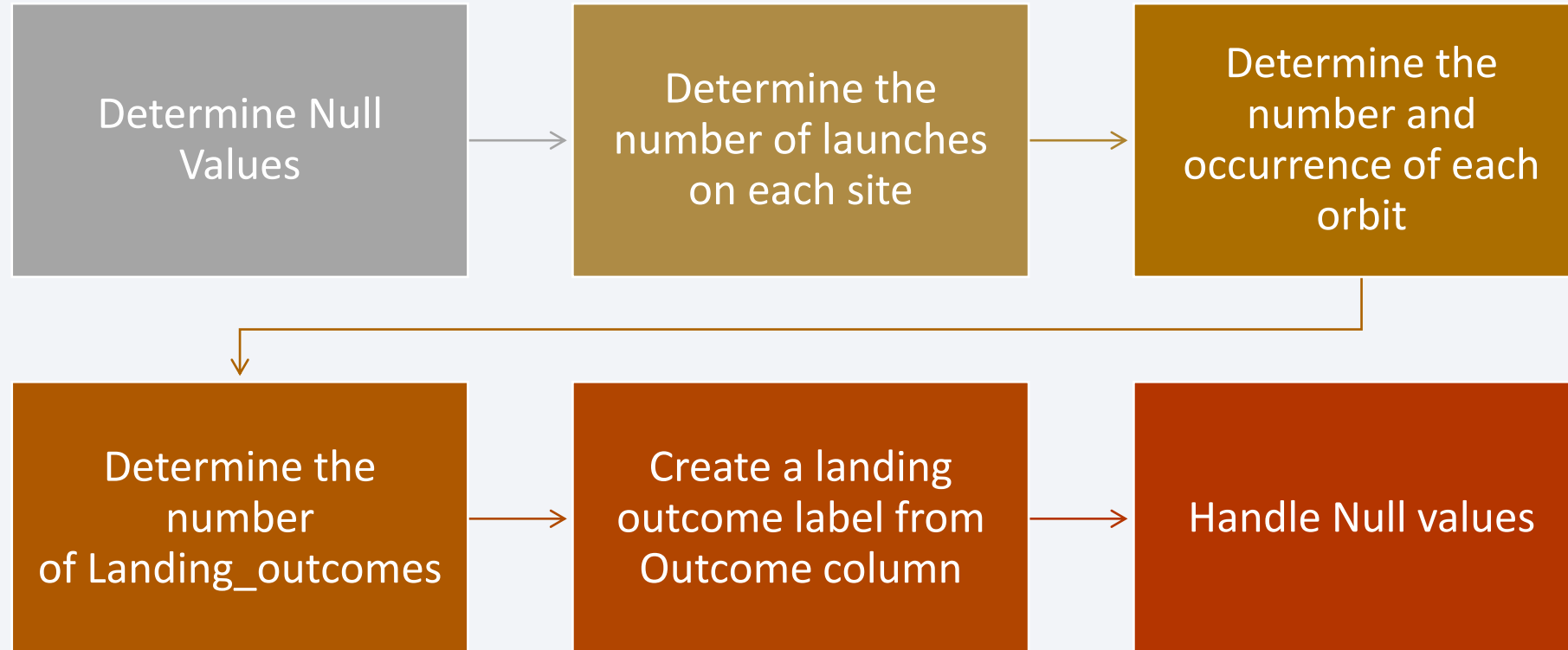
In [18]:

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 121 entries, 0 to 120
Data columns (total 11 columns):
 #   Column              Non-Null Count  Dtype
---  -
 0   Flight No.          121 non-null   object
 1   Launch site         121 non-null   object
 2   Payload             121 non-null   object
 3   Payload mass        121 non-null   object
 4   Orbit               121 non-null   object
 5   Customer            120 non-null   object
 6   Launch outcome      121 non-null   object
 7   Version Booster     121 non-null   object
 8   Booster landing     121 non-null   object
 9   Date               121 non-null   object
10   Time               121 non-null   object
dtypes: object(11)
memory usage: 10.5+ KB
```

<https://github.com/Mamad66Ahmadi/IBM-Data-Science-Capstone/blob/Module1/jupyter-labs-webscraping.ipynb>

Data Wrangling



https://github.com/Mamad66Ahmadi/IBM-Data-Science-Capstone/blob/Module1/labs-jupyter-spacex-data_wrangling_jupyterlite.jupyterlite.ipynb

EDA with Data Visualization

- For machine learning models we should get insight into the relationship between parameters. Therefore, various type of plots such as: Scatter plots, line charts, and bar plots are used to show the relation between different parameters.

Flight Number vs. Payload Mass

Flight Number vs. Launch Site

Payload Mass vs. Launch Site

Orbit vs. Success Rate

Flight Number vs. Orbit

Payload vs Orbit

Success Yearly Trend

<https://github.com/Mamad66Ahmadi/IBM-Data-Science-Capstone/blob/Module2/jupyter-labs-eda-dataviz.ipynb.jupyterlite.ipynb>

EDA with SQL

- Displayed the names of the unique launch sites in the space mission
- Displayed 5 records where launch sites begin with the string 'CCA'
- Displayed the total payload mass carried by boosters launched by NASA (CRS)
- Displayed average payload mass carried by booster version F9 v1.1
- Listed the date when the first succesful landing outcome in ground pad was acheived.
- Listed the names of the boosters which have success in drone ship and have payload mass greater than 4000 but less than 6000
- Listed the total number of successful and failure mission outcomes
- Listed the names of the booster_versions which have carried the maximum payload mass. Use a subquery
- Listed the records which will display the month names, failure landing_outcomes in drone ship ,booster versions, launch_site for the months in year 2015.
- Ranked the count of successful landing_outcomes between the date 04-06-2010 and 20-03-2017 in descending order.

https://github.com/Mamad66Ahmadi/IBM-Data-Science-Capstone/blob/Module2/jupyter-labs-eda-sql-coursera_sqllite.ipynb

Build an Interactive Map with Folium

with aim to finding an optimal location for building a launch site:

- Calculated the distances between a launch site to its proximities
- Marked down a point on the closest coastline using MousePosition and calculated the distance between the coastline point and the launch site
- Drew a PolyLine between a launch site to the selected coastline point
- Created a marker with distance to a closest city, railway, highway, etc.
- Drew a line between the marker to the launch site

https://github.com/Mamad66Ahmadi/IBM-Data-Science-Capstone/blob/Module3/lab_jupyter_launch_site_location.jupyterlite.ipynb

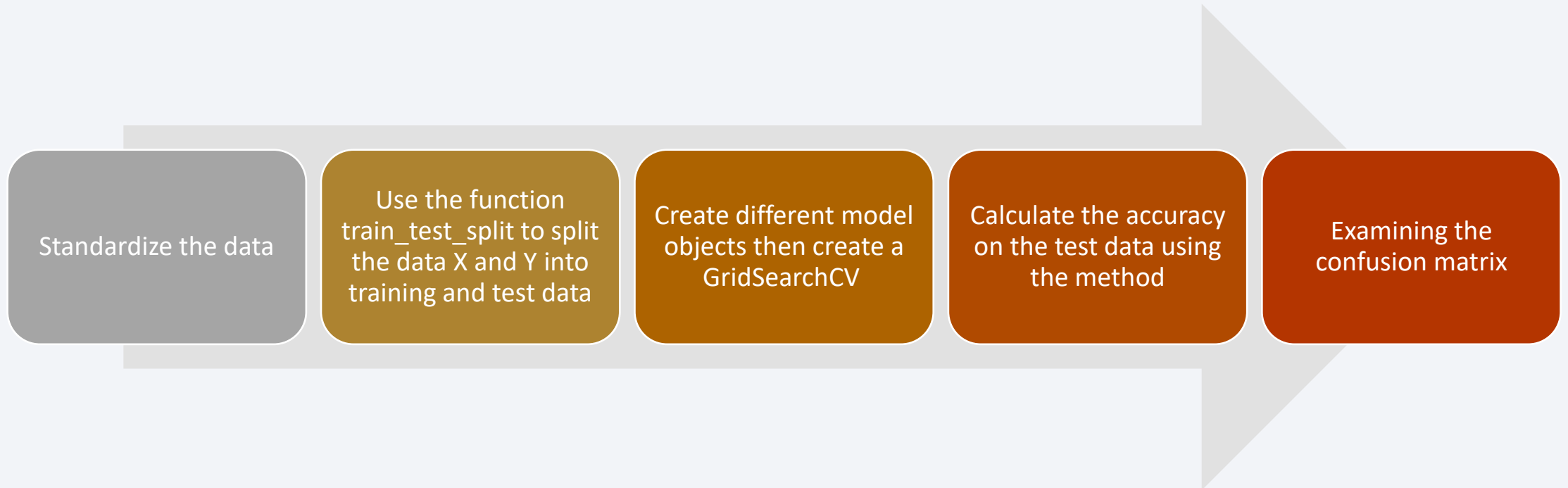
Build a Dashboard with Plotly Dash

Dashboard with a selectable pie chart and a scatter plot.

- Pie chart shows distribution of successful landings across all launch sites and can be selected to show individual launch site success rates.
- The pie chart is used to visualize launch site success rate.
- Scatter plot takes two inputs: All sites or individual site and payload mass on a slider between 0 and 10000 kg.
- The scatter plot can help us see how success varies across launch sites, payload mass, and booster version category.

https://github.com/Mamad66Ahmadi/IBM-Data-Science-Capstone/blob/Module3/spacex_dash_app.py

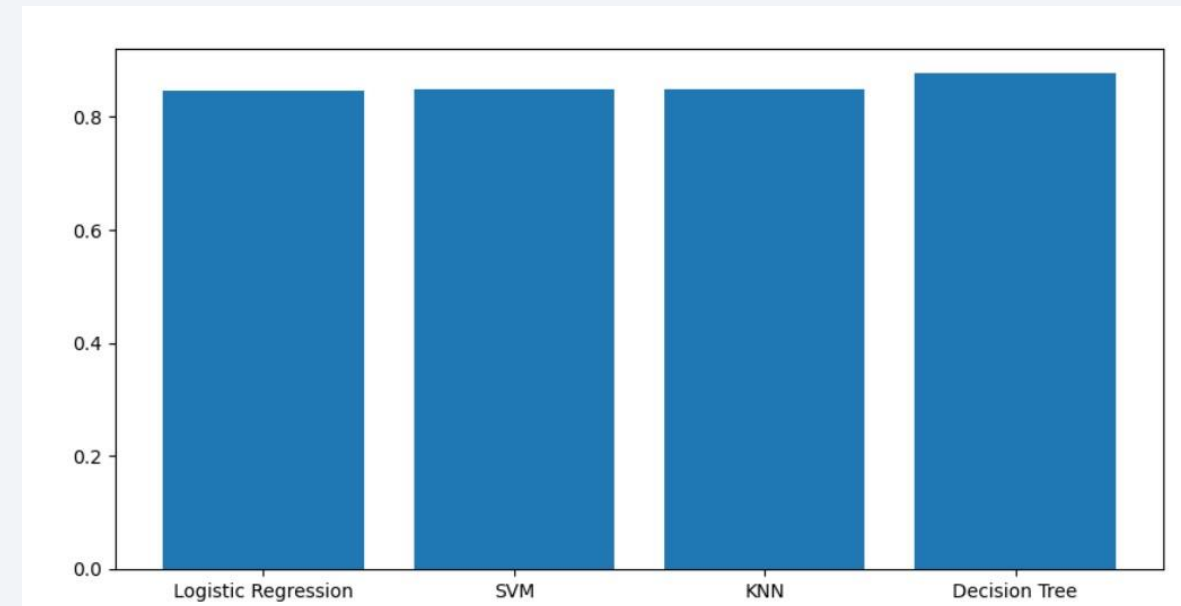
Predictive Analysis (Classification)



https://github.com/Mamad66Ahmadi/IBM-Data-Science-Capstone/blob/Module4/SpaceX_Machine_Learning_Prediction_Part_5.jupyterlite.ipynb

Results

- The SVM, KNN, and Logistic Regression models are the best in terms of prediction accuracy for this dataset.
- Low weighted payloads perform better than the heavier payloads.
- The success rates for SpaceX launches is directly proportional time in years they will eventually perfect the launches.
- KSC LC 39A had the most successful launches from all the sites.
- Orbit GEO,HEO,SSO,ES L1 has the best Success Rate.



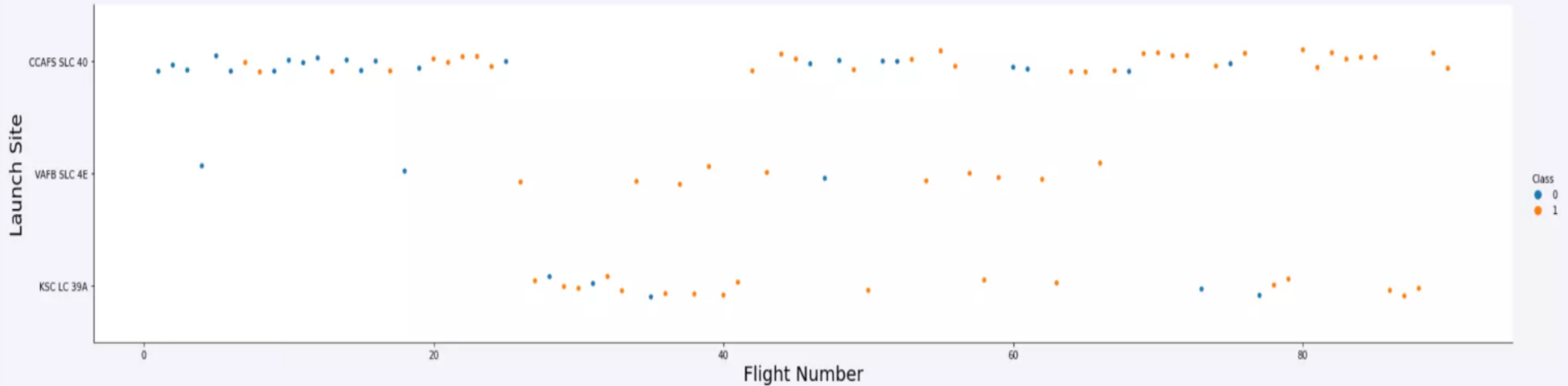
The background of the slide is an abstract composition. It features a dark blue base color. Overlaid on this are numerous diagonal streaks in shades of red and cyan. A faint, light blue grid pattern is also visible, particularly in the lower half of the image. The overall effect is dynamic and technological.

Section 2

Insights drawn from EDA

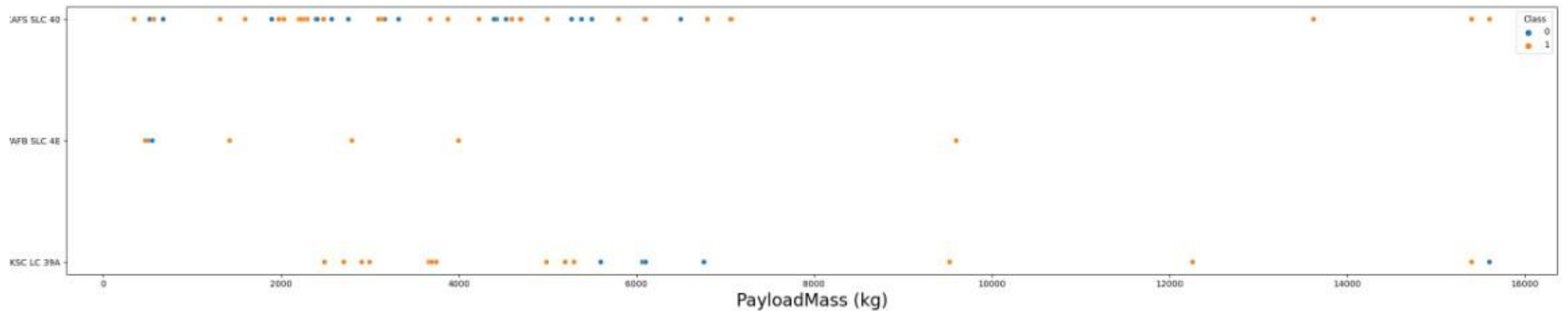
Flight Number vs. Launch Site

- Launches from the site of CCAFS SLC 40 are significantly higher than lunches from other sites



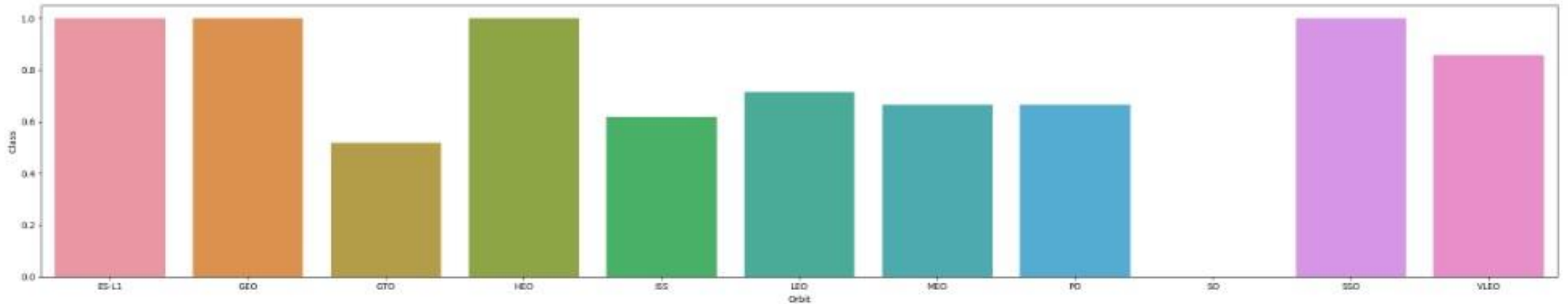
Payload vs. Launch Site

- The majority of Pay Loads with lower Mass have been launched from CCAFS SLC 40.



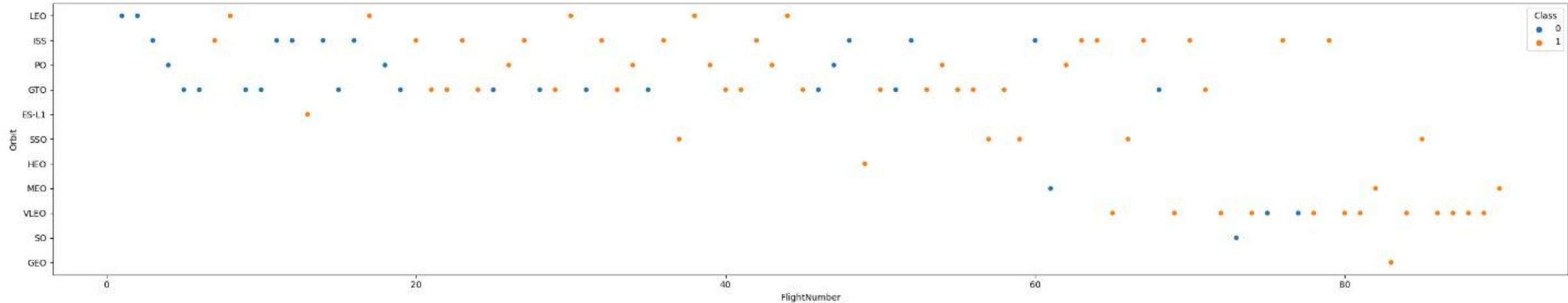
Success Rate vs. Orbit Type

- ESL1, GEO, HEO and SSO orbits had the highest success rate.



Flight Number vs. Orbit Type

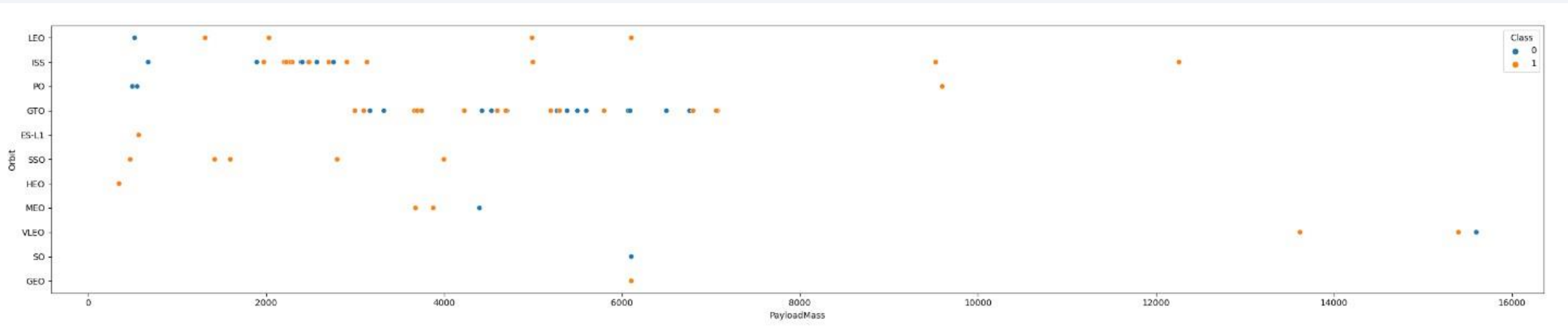
- There is an increasing rate of VLEO launches in recent years.



Payload vs. Orbit Type

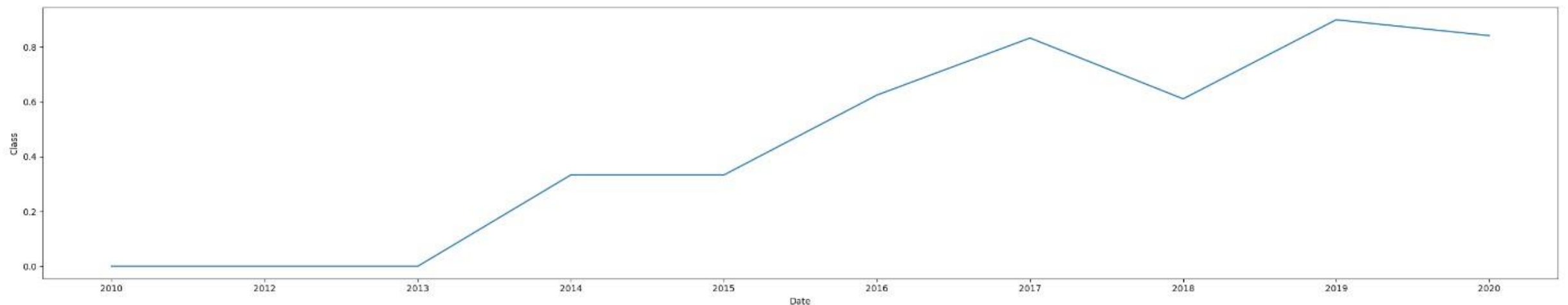
At ISS most of payloads were around 2000

At GTO most of payload were between 3000 and 7000



Launch Success Yearly Trend

- Success rate has increased incredibly in recent years



All Launch Site Names

In [7]: `%sql select distinct launch_site from SPACEXTBL`

`* sqlite:///my_data1.db`

Done.

Out[7]: **Launch_Site**

CCAFS LC-40

VAFB SLC-4E

KSC LC-39A

CCAFS SLC-40

None

Launch Site Names Begin with 'CCA'

```
In [14]: %sql select * from SPACEXTBL where launch_site like 'CCA%' limit 5
```

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

```
Done.
```

```
Out[14]:
```

Date	Time (UTC)	Booster_Version	Launch_Site	Payload	PAYLOAD_MASS_KG	Orbit	Customer	Mission_Outcome	Landing_Outc
06/04/2010	18:45:00	F9 v1.0 B0003	CCAFS LC-40	Dragon Spacecraft Qualification Unit	0.0	LEO	SpaceX	Success	Failure (parachute)
12/08/2010	15:43:00	F9 v1.0 B0004	CCAFS LC-40	Dragon demo flight C1, two CubeSats, barrel of Brouere cheese	0.0	LEO (ISS)	NASA (COTS) NRO	Success	Failure (parachute)
22/05/2012	7:44:00	F9 v1.0 B0005	CCAFS LC-40	Dragon demo flight C2	525.0	LEO (ISS)	NASA (COTS)	Success	No attempt
10/08/2012	0:35:00	F9 v1.0 B0006	CCAFS LC-40	SpaceX CRS-1	500.0	LEO (ISS)	NASA (CRS)	Success	No attempt
03/01/2013	15:10:00	F9 v1.0 B0007	CCAFS LC-40	SpaceX CRS-2	677.0	LEO (ISS)	NASA (CRS)	Success	No attempt

Total Payload Mass

```
In [9]: %sql select sum(PAYLOAD_MASS_KG_) from SPACEXTBL where customer = 'NASA (CRS)'
```

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

```
Out[9]: sum(PAYLOAD_MASS_KG_)  
         45596.0
```

Average Payload Mass by F9 v1.1

```
In [10]: %sql select avg(PAYLOAD_MASS_KG_) from SPACEXTBL where Booster_Version = 'F9 v1.1'
```

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

```
Out[10]: avg(PAYLOAD_MASS_KG_)  
                2928.4
```

First Successful Ground Landing Date

```
In [24]: %sql select min(Date) from SPACEXTBL where "Landing_Outcome" = "Success (ground pad)"
* sqlite:///my_data1.db
Done.
Out[24]: min(Date)
          01/08/2018
```


Successful Drone Ship Landing with Payload between 4000 and 6000

%sql select booster_version from SPACEXTBL where "Landing_Outcome"='Success (drone ship)' and PAYLOAD_MASS__KG_ between 4000 and 6000

```
In [26]: %sql select booster_version from SPACEXTBL where "Landing_Outcome"='Success (drone ship)' and PAYLOAD_MASS__KG_ between 4000 and 6000
```

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

```
Out[26]: Booster_Version
```

```
F9 FT B1022
```

```
F9 FT B1026
```

```
F9 FT B1021.2
```

```
F9 FT B1031.2
```

Total Number of Successful and Failure Mission Outcomes

In [27]: `%sql select Mission_Outcome, count(*) from SPACEXTBL group by Mission_Outcome`

`* sqlite:///my_data1.db`

Done.

Out[27]:

Mission_Outcome	count(*)
None	898
Failure (in flight)	1
Success	98
Success	1
Success (payload status unclear)	1

Boosters Carried Maximum Payload

```
In [28]: %sql select distinct Booster_Version from SPACEXTBL where PAYLOAD_MASS__KG_ = (select max(PAYLOAD_MASS__KG_) from SPACEXTBL
```

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

```
Done.
```

```
Out[28]: 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

%sql select substr(Date, 4, 2) as month, "Landing_Outcome", Booster_Version , Launch_Site from SPACEXTBL where "Landing_Outcome" = 'Failure (drone ship)' and substr(Date,7,4)='2015'

In [31]: %sql select substr(Date, 4, 2) as month, "Landing_Outcome", Booster_Version , Launch_Site from SPACEXTBL where "Landing_Out

* sqlite:///my_data1.db

Done.

Out[31]:

month	Landing_Outcome	Booster_Version	Launch_Site
10	Failure (drone ship)	F9 v1.1 B1012	CCAFS LC-40
04	Failure (drone ship)	F9 v1.1 B1015	CCAFS LC-40

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

- %sql SELECT "Landing_Outcome", count(*) AS count, RANK() OVER (ORDER BY count(*) DESC) AS rank FROM SPACEXTBL WHERE "Landing_Outcome" like 'Success%' and Date between '04-06-2010' and '20-03-2017' GROUP BY "Landing_Outcome" ORDER BY count DESC

In [32]: %sql SELECT "Landing_Outcome", count(*) AS count, RANK() OVER (ORDER BY count(*) DESC) AS rank FROM SPACEXTBL WHERE "Landing

* sqlite:///my_data1.db

Done.

Out[32]:

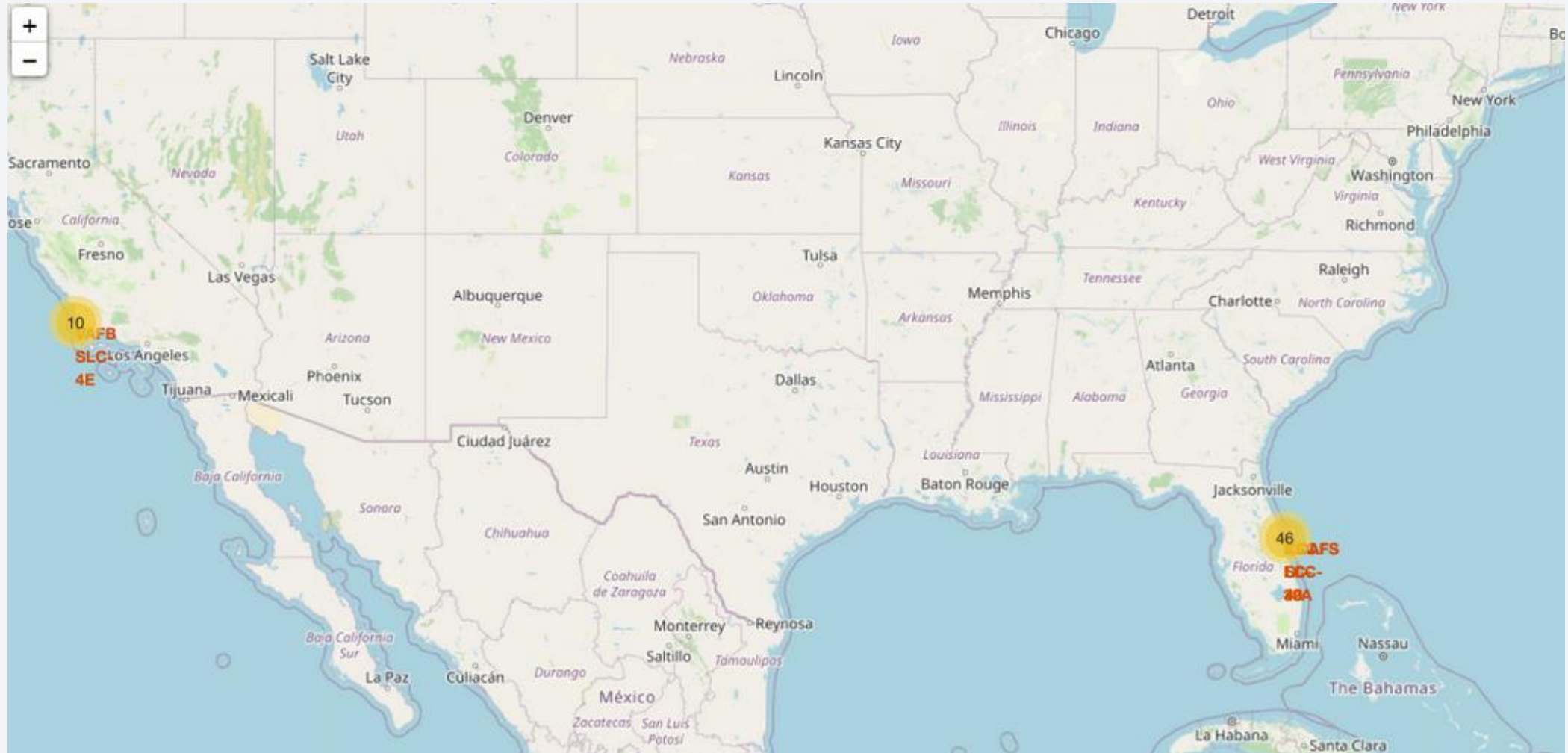
Landing_Outcome	count	rank
Success	20	1
Success (drone ship)	8	2
Success (ground pad)	7	3

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.

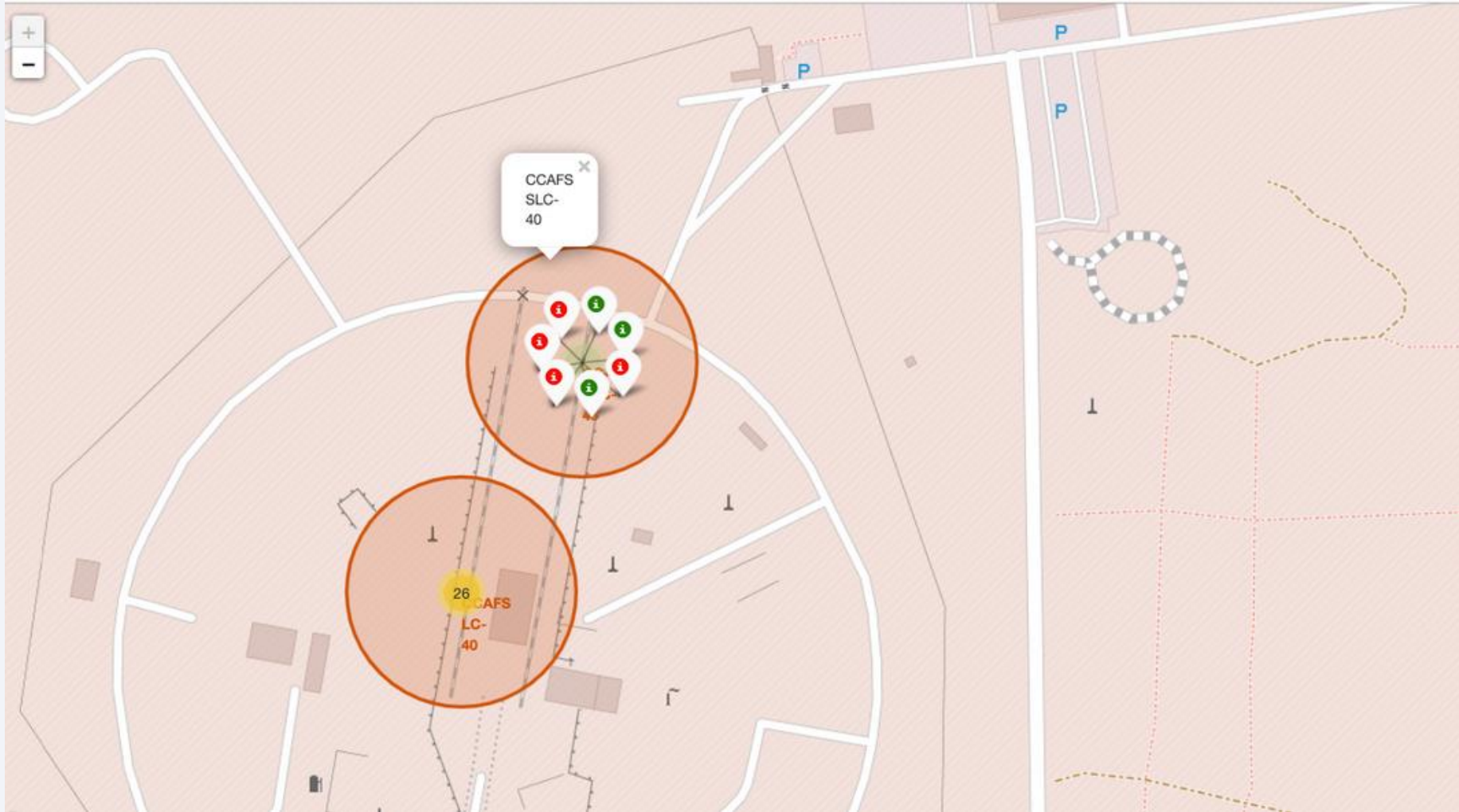
Section 3

Launch Sites Proximities Analysis

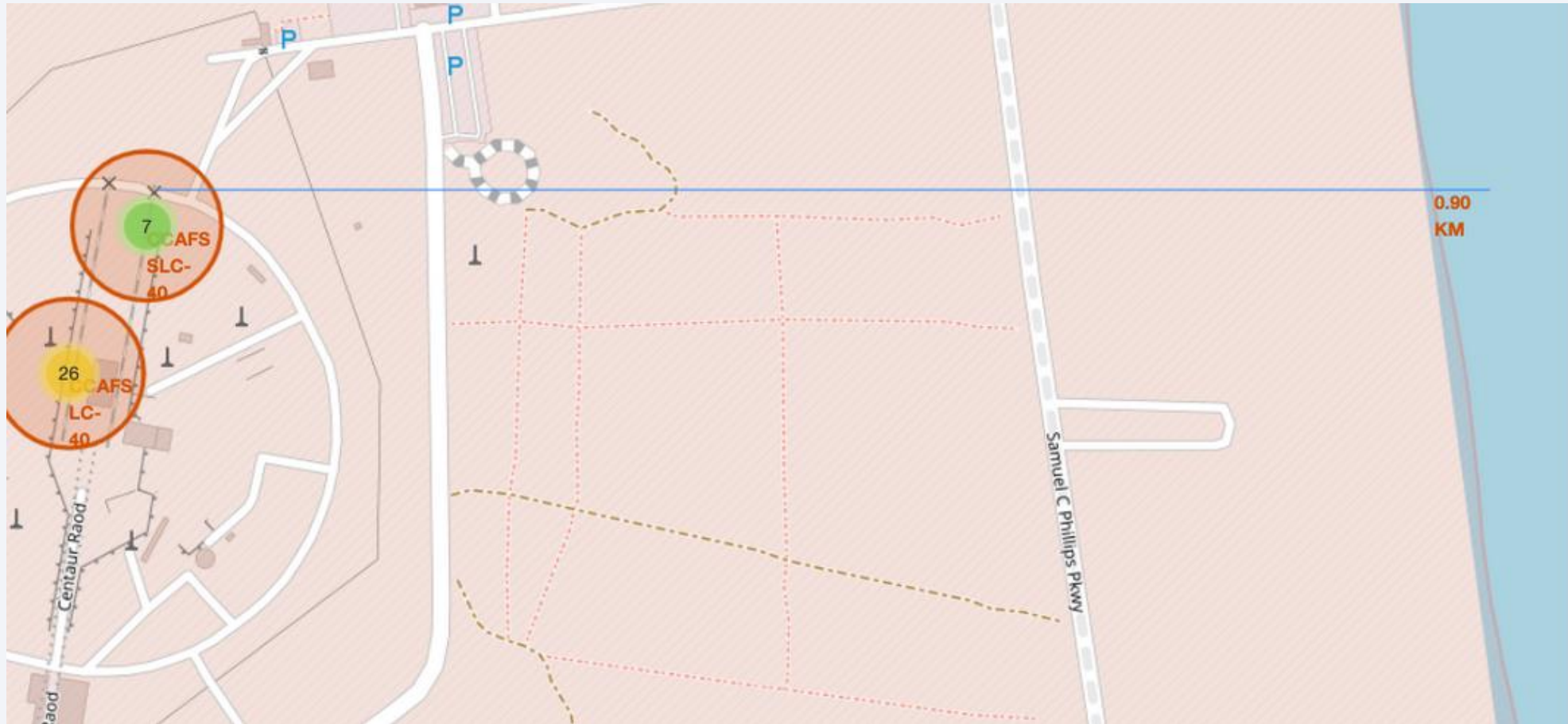
<Folium Map Screenshot 1>



<Folium Map Screenshot 2>



<Folium Map Screenshot 3>

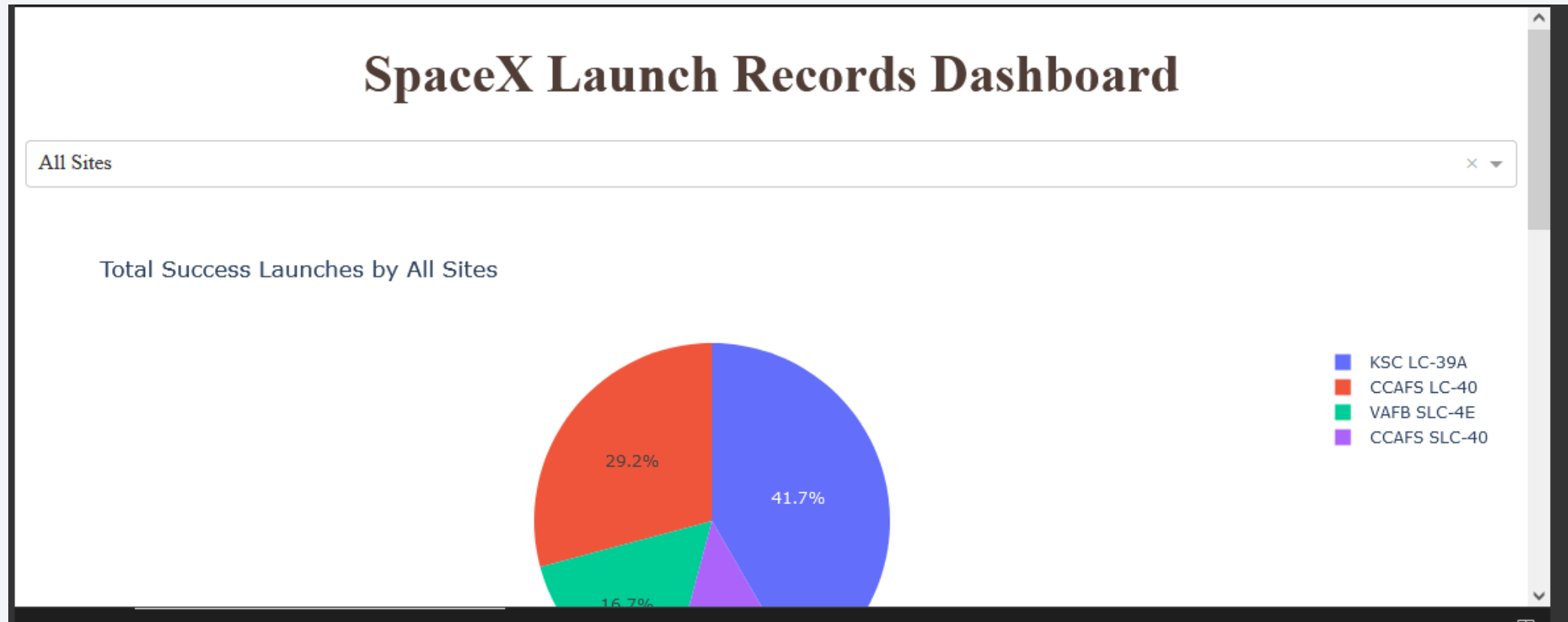




Section 4

Build a Dashboard with Plotly Dash

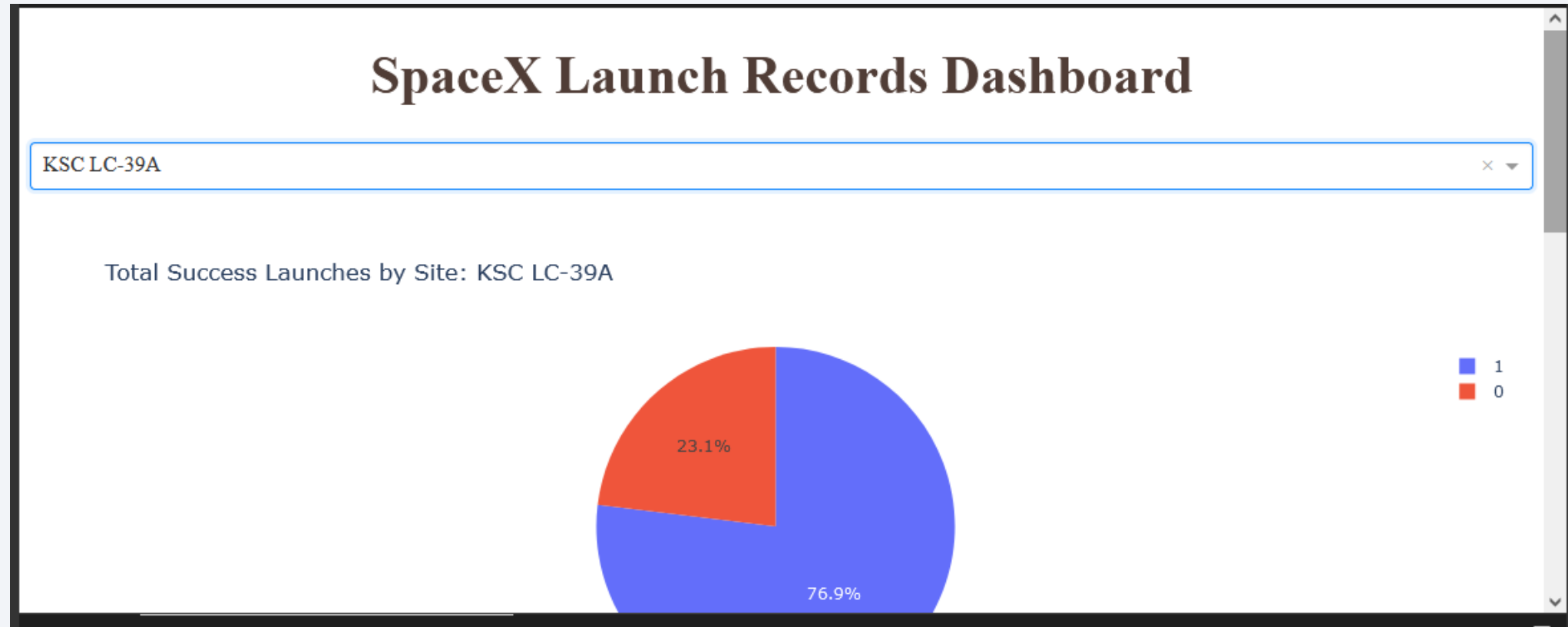
<Dashboard Screenshot 1>



<Dashboard Screenshot 2>



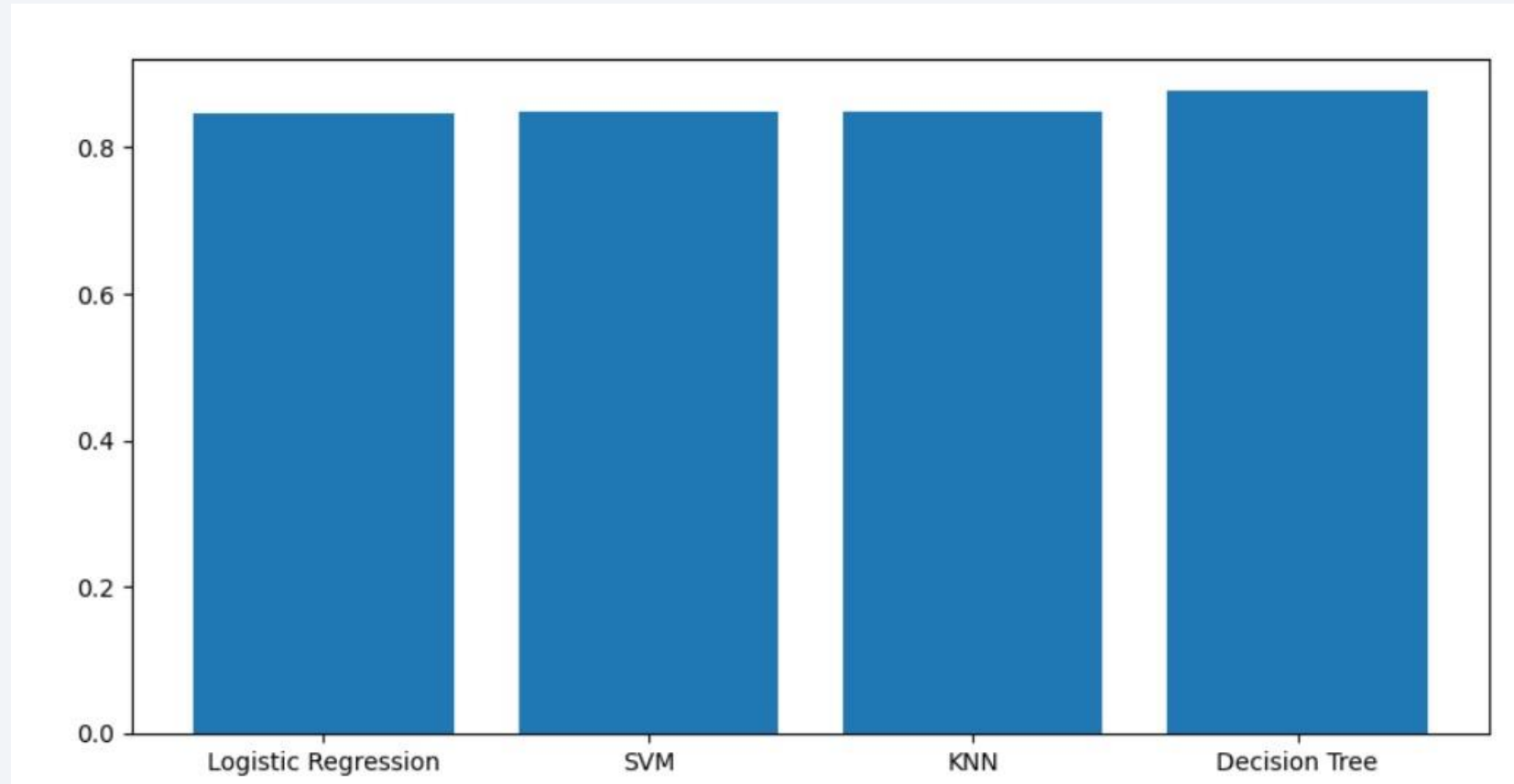
<Dashboard Screenshot 3>



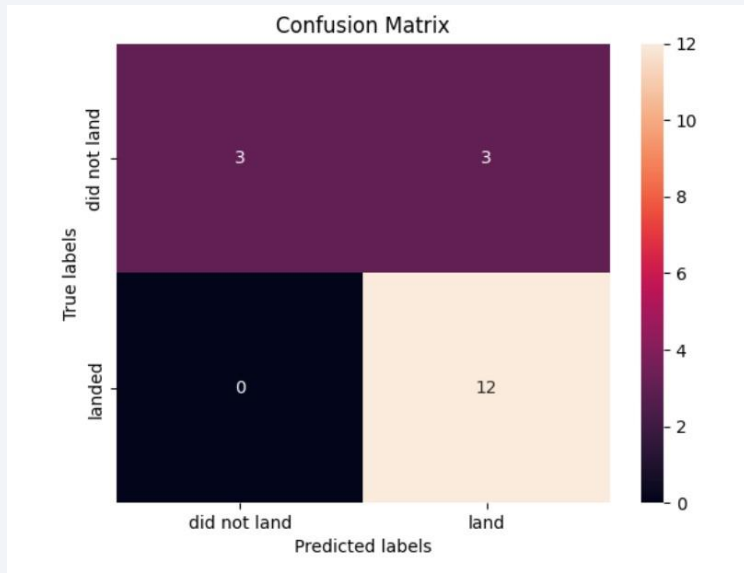
Section 5

Predictive Analysis (Classification)

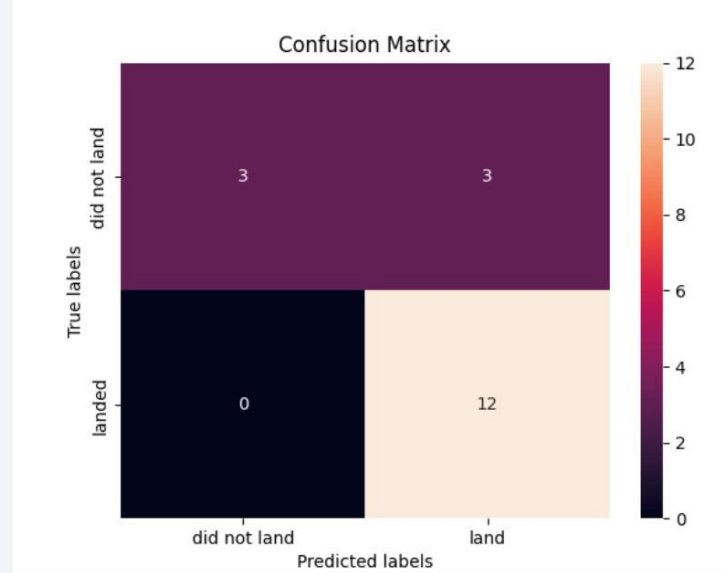
Classification Accuracy



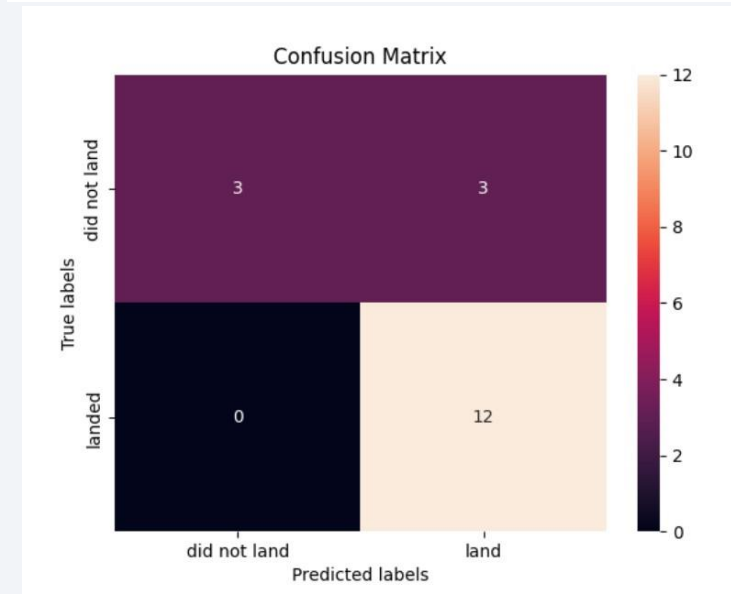
Confusion Matrix



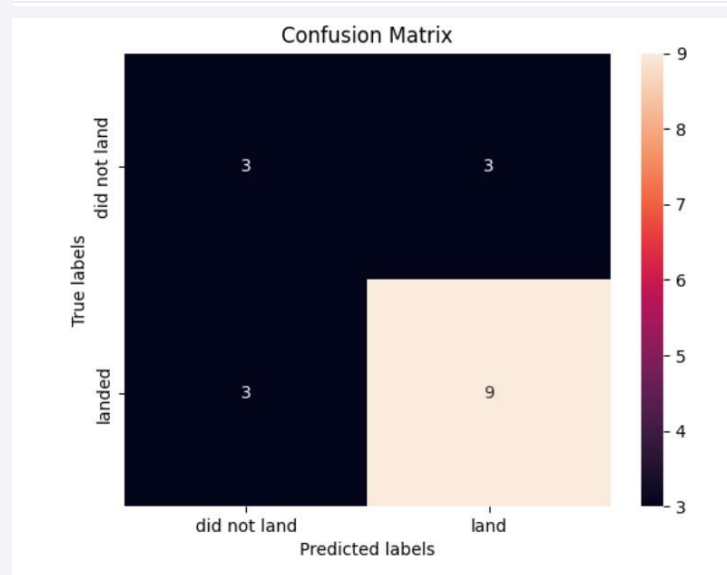
Logistic Regression



SVM



KNN



Decision Tree

Conclusions

- The Decision Tree model has the best prediction accuracy.
- Low weighted payloads has a better success rate than the heavier payloads.
- The success rates for SpaceX launches have increased significantly in recent years.
- KSC LC 39A had the most successful launches among all other sites.
- Orbit GEO,HEO,SSO,ES L1 had the best Success Rate.

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!

