



IBM Developer
SKILLS NETWORK

Winning Space Race with Data Science

Ishanee Neb
19 June 2025



Outline

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

Executive Summary

- In this report, we predict whether the SpaceX Falcon 9 first stage will land successfully.
- The main steps in this project include:
 - Data collection, wrangling, and formatting
 - Exploratory data analysis
 - Interactive data visualization
 - Machine learning prediction
- In this report, it is also concluded that decision tree may be the best machine learning algorithm to predict if the Falcon 9 first stage will land successfully.

Introduction

- Through this report, we can predict if the Falcon 9 first stage will land successfully. SpaceX advertises Falcon 9 rocket launches on its website with a cost of 62 million dollars; other providers cost upward of 165 million dollars each, much of the savings is because SpaceX can reuse the first stage. Therefore if we can determine if the first stage will land, we can determine the cost of a launch.
- The main question that we are trying to answer is, for a given set of features about a Falcon 9 rocket launch which include its payload mass, orbit type, launch site, and so on, will the first stage of the rocket land successfully?

Section 1

Methodology

Methodology

Executive Summary

- Data collection methodology:
 - Data collection, wrangling, and formatting, using: a) SpaceX API b) Web scraping
- Perform exploratory data analysis (EDA) using visualization and SQL
 - Exploratory data analysis (EDA), using: a) Pandas and NumPy b) SQL
- Perform interactive visual analytics
 - Data visualization, using: a) Matplotlib and Seaborn b) Folium c) Dash
- Perform predictive analysis using classification models
 - Machine learning prediction, using a) Logistic regression b) Support vector machine (SVM) c) Decision tree d) K-nearest neighbors (KNN)

Data collection methodology

• SpaceX API

- The API used is <https://api.spacexdata.com/v4/rockets/>.
- The API provides data about many types of rocket launches done by SpaceX, the data is therefore filtered to include only Falcon 9 launches.
- Every missing value in the data is replaced the mean the column that the missing value belongs to.
- We end up with 90 rows or instances and 17 columns or features. The picture below shows the first few rows of the data:

FlightNumber	Date	BoosterVersion	PayloadMass	Orbit	LaunchSite	Outcome	Flights	GridFins	Reused	Legs	LandingPad	Block	ReusedCount	Serial	Longitude	Latitude	
4	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
5	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
6	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
7	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
8	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

• Web scraping

- The data is scraped from https://en.wikipedia.org/w/index.php?title=List_of_Falcon_9_and_Falcon_Heavy_launches&oldid=1027686922
- The website contains only the data about Falcon 9 launches.
- We end up with 121 rows or instances and 11 columns or features. The picture below shows the first few rows of the data:

	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

Exploratory data analysis (EDA)

- **Pandas and NumPy**

- Functions from the Pandas and NumPy libraries are used to derive basic information about the data collected, which includes:
 - The number of launches on each launch site
 - The number of occurrences of each orbit
 - The number and occurrence of each mission outcome

- **SQL**

- The data is queried using SQL to answer several questions about the data such as:
 - The names of the unique launch sites in the space mission
 - The total payload mass carried by boosters launched by NASA (CRS)
 - The average payload mass carried by booster version F9 v1.1

Data Visualization

• Matplotlib and Seaborn

- Functions from the Matplotlib and Seaborn libraries are used to visualize the data through scatterplots, bar charts, and line charts.
- The plots and charts are used to understand more about the relationships between several features, such as:
 - The relationship between flight number and launch site
 - The relationship between payload mass and launch site
 - The relationship between success rate and orbit type

• Folium

- Functions from the Folium libraries are used to visualize the data through interactive maps.
- The Folium library is used to:
 - Mark all launch sites on a map
 - Mark the succeeded launches and failed launches for each site on the map
 - Mark the distances between a launch site to its proximities such as the nearest city, railway, or highway

• Dash

- Functions from Dash are used to generate an interactive site where we can toggle the input using a dropdown menu and a range slider.
- Using a pie chart and a scatterplot, the interactive site shows:
 - The total success launches from each launch site
 - The correlation between payload mass and mission outcome (success or failure) for each launch site

Machine Learning Prediction

- Functions from the Scikit-learn library are used to create our machine learning models.
- The machine learning prediction phase include the following steps:
 - Standardizing the data
 - Splitting the data into training and test data
 - Creating machine learning models, which include:
 - Logistic regression
 - Support vector machine (SVM)
 - Decision tree
 - K nearest neighbors (KNN)
 - Fit the models on the training set
 - Find the best combination of hyperparameters for each model
 - Evaluate the models based on their accuracy scores and confusion matrix

Results

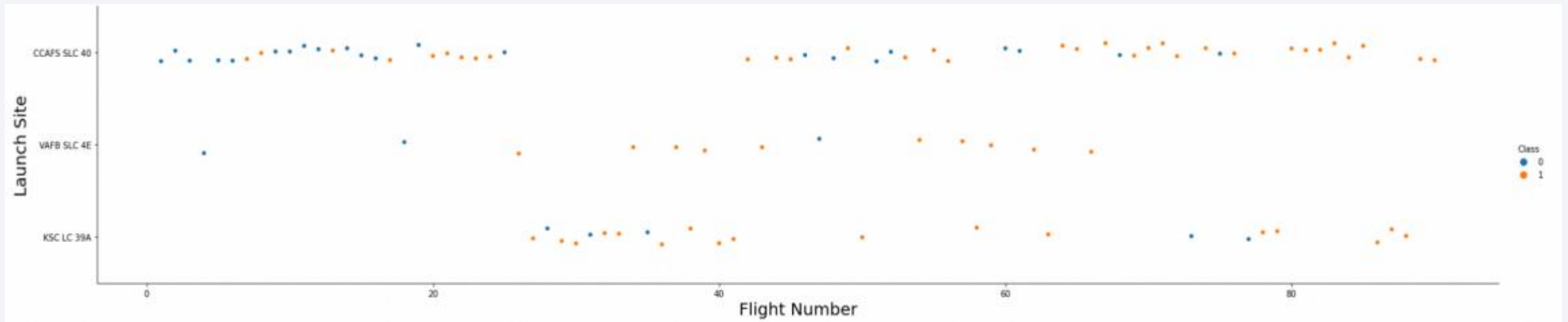
- The results are split into 5 sections:
 - SQL (EDA with SQL)
 - Matplotlib and Seaborn (EDA with Visualization)
 - Folium
 - Dash
 - Predictive Analysis
- In all of the graphs that follow, class 0 represents a failed launch outcome while class 1 represents a successful launch outcome

The background of the slide is an abstract composition. It features a dark blue field on the left side, which transitions into a complex pattern of diagonal streaks and lines in shades of blue, red, and cyan on the right. These streaks have a textured, almost woven appearance, suggesting a digital or data-driven theme. The overall effect is dynamic and modern.

Section 2

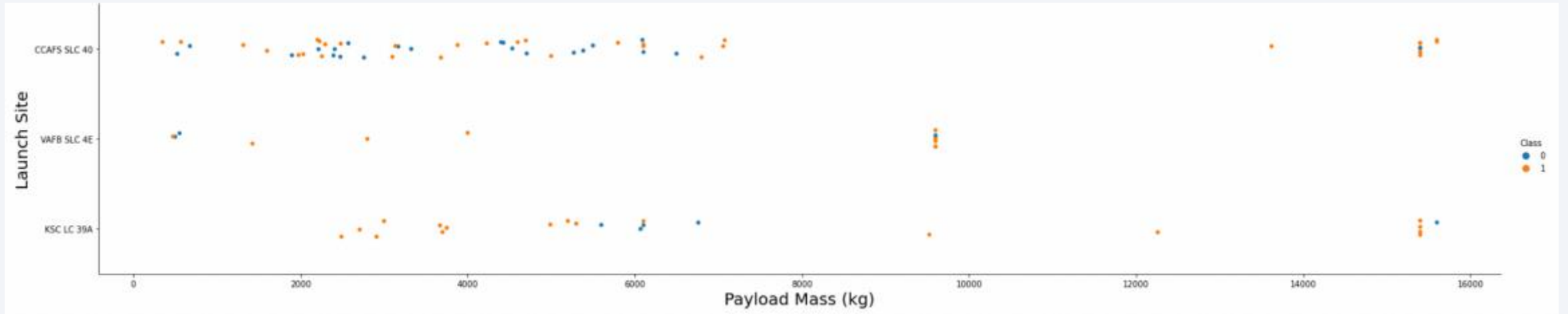
Insights drawn from EDA

Flight Number vs. Launch Site



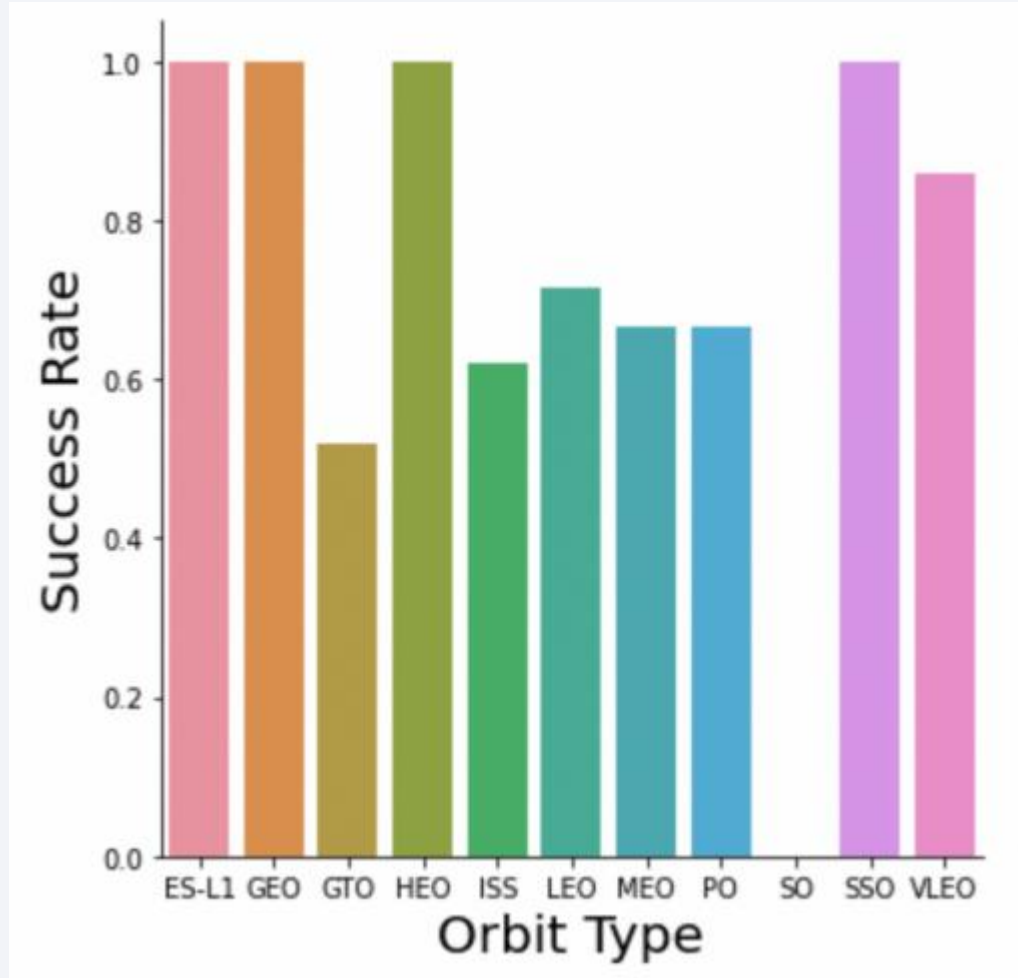
- Explanation:
 - The earliest flights all failed while the latest flights all succeeded.
 - The CCAFS SLC 40 launch site has about a half of all launches.
 - VAFB SLC 4E and KSC LC 39A have higher success rates.
 - It can be assumed that each new launch has a higher rate of success.

Payload vs. Launch Site



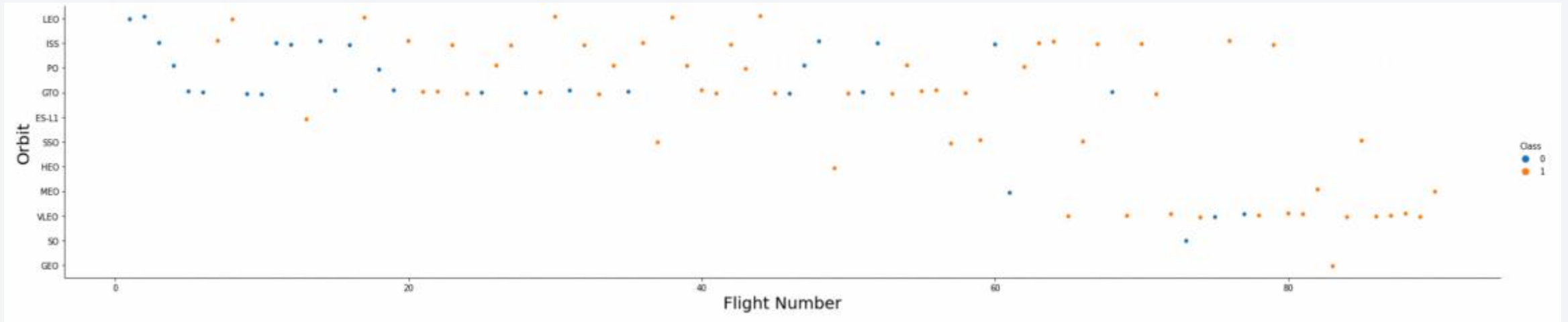
- Explanation:
 - For every launch site the higher the payload mass, the higher the success rate.
 - Most of the launches with payload mass over 7000 kg were successful.
 - KSC LC 39A has a 100% success rate for payload mass under 5500 kg too.

Success Rate vs. Orbit Type



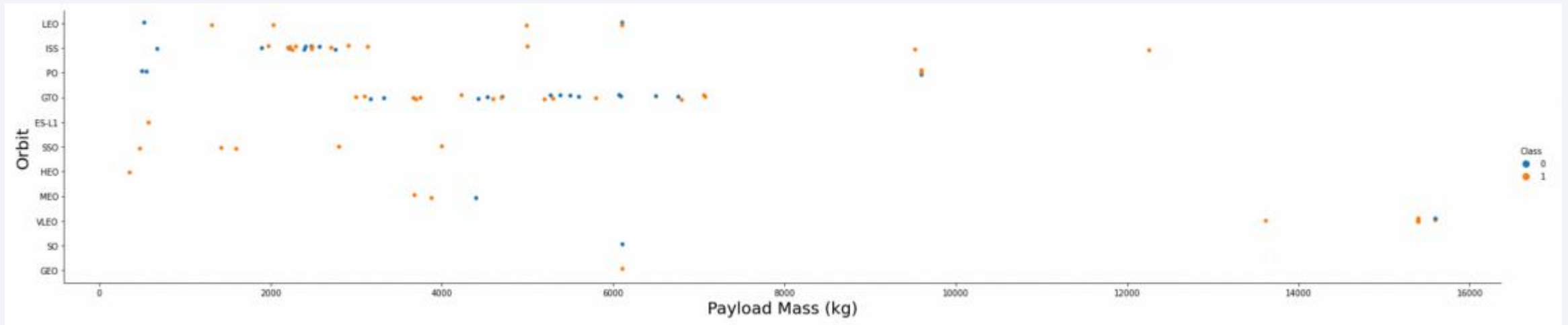
- Explanation:
 - Orbits with 100% success rate: - ES-L1, GEO, HEO, SSO
 - Orbits with 0% success rate: - SO
 - Orbits with success rate between 50% and 85%: - GTO, ISS, LEO, MEO, PO

Flight Number vs. Orbit Type



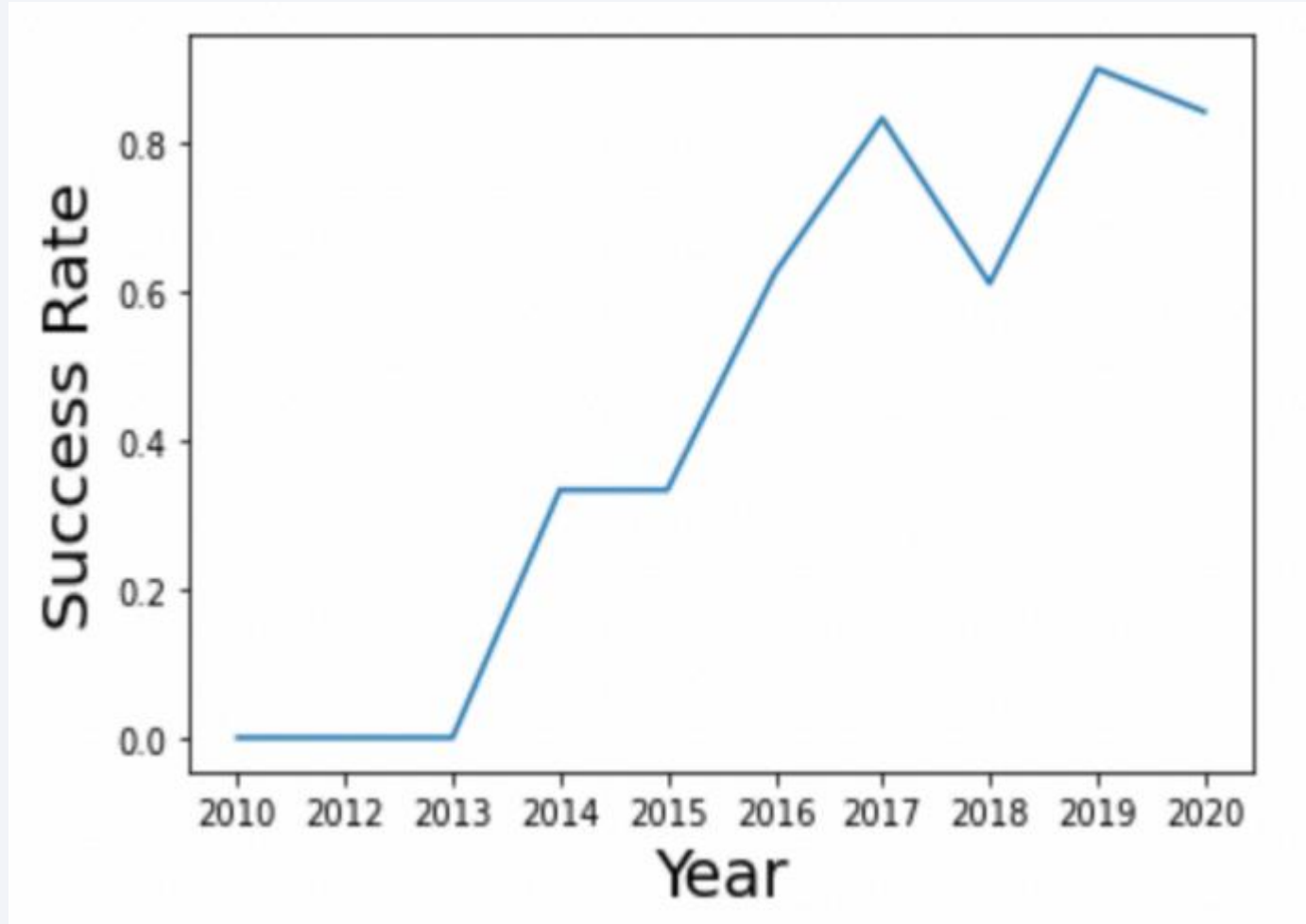
- Explanation:
 - In the LEO orbit the Success appears related to the number of flights; on the other hand, there seems to be no relationship between flight number when in GTO orbit.

Payload vs. Orbit Type



- Explanation:
 - Heavy payloads have a negative influence on GTO orbits and positive on GTO and Polar LEO (ISS) orbits.

Launch Success Yearly Trend



- Explanation:
 - The success rate since 2013 kept increasing till 2020.

All Launch Site Names

```
In [11]: %sql select DISTINCT "Launch_Site" from SPACEXTABLE
```

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

```
Out[11]: Launch_Site  
-----  
          CCAFS LC-40  
          VAFB SLC-4E  
          KSC LC-39A  
          CCAFS SLC-40
```

- Explanation:
 - Displaying the names of the unique launch sites in the space mission.

Launch Site Names Begin with 'CCA'

In [12]: `%sql SELECT * FROM SPACEXTABLE WHERE "Launch_Site" LIKE 'CCA%' LIMIT 5`

* sqlite:///my_data1.db
Done.

Out[12]:

Date	Time (UTC)	Booster_Version	Launch_Site	Payload	PAYLOAD_MASS_KG_	Orbit	Customer	Mission_Outcome	Landing_Outcome
2010-06-04	18:45:00	F9 v1.0 B0003	CCAFS LC-40	Dragon Spacecraft Qualification Unit	0	LEO	SpaceX	Success	Failure (parachute)
2010-12-08	15:43:00	F9 v1.0 B0004	CCAFS LC-40	Dragon demo flight C1, two CubeSats, barrel of Brouere cheese	0	LEO (ISS)	NASA (COTS) NRO	Success	Failure (parachute)
2012-05-22	7:44:00	F9 v1.0 B0005	CCAFS LC-40	Dragon demo flight C2	525	LEO (ISS)	NASA (COTS)	Success	No attempt
2012-10-08	0:35:00	F9 v1.0 B0006	CCAFS LC-40	SpaceX CRS-1	500	LEO (ISS)	NASA (CRS)	Success	No attempt
2013-03-01	15:10:00	F9 v1.0 B0007	CCAFS LC-40	SpaceX CRS-2	677	LEO (ISS)	NASA (CRS)	Success	No attempt

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

Total Payload Mass

```
In [13]: %sql SELECT SUM("PAYLOAD_MASS_KG_") FROM SPACEXTABLE WHERE "Customer"="NASA (CRS)"
* sqlite:///my_data1.db
Done.
Out[13]: SUM("PAYLOAD_MASS_KG_")
         45596
```

- Explanation:
 - Displaying the total payload mass carried by boosters launched by NASA (CRS).

Average Payload Mass by F9 v1.1

```
In [14]: %sql SELECT AVG(PAYLOAD_MASS_KG_) FROM SPACEXTABLE WHERE "Booster_Version" LIKE "F9 v1.1%"
* sqlite:///my_data1.db
Done.
Out[14]: AVG(PAYLOAD_MASS_KG_)
          2534.6666666666665
```

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

First Successful Ground Landing Date

```
In [16]: %sql SELECT MIN("Date") FROM SPACEXTABLE WHERE "Landing_Outcome"="Success (ground pad)"
* sqlite:///my_data1.db
Done.
Out[16]: MIN("Date")
         2015-12-22
```

- Explanation:
 - Listing the date when the first successful landing outcome in ground pad was achieved.

Successful Drone Ship Landing with Payload between 4000 and 6000

```
In [18]: %sql SELECT "Booster_Version" FROM SPACEXTABLE WHERE "Landing_Outcome"="Success (drone ship)" AND "PAYLOAD_MASS__KG_">4000 AND "PAYLOAD_MASS__KG_"<6000
```

* sqlite:///my_data1.db
Done.

Out[18]: **Booster_Version**

F9 FT B1022
F9 FT B1026
F9 FT B1021.2
F9 FT B1031.2

- Explanation: Listing the names of the boosters which have success in drone ship and have payload mass greater than 4000 but less than 6000.

Total Number of Successful and Failure Mission Outcomes

```
In [11]: %sql SELECT "Landing_Outcome",COUNT(*) FROM SPACEXTABLE GROUP BY "Landing_Outcome"
```

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

```
Out[11]:
```

Landing_Outcome	COUNT(*)
Controlled (ocean)	5
Failure	3
Failure (drone ship)	5
Failure (parachute)	2
No attempt	21
No attempt	1
Precluded (drone ship)	1
Success	38
Success (drone ship)	14
Success (ground pad)	9
Uncontrolled (ocean)	2

- Explanation:
 - Listing the total number of successful and failure mission outcomes.

Boosters Carried Maximum Payload

```
In [22]: %sql SELECT "Booster_Version" FROM SPACEXTABLE WHERE "PAYLOAD_MASS__KG_"= (SELECT MAX("PAYLOAD_MASS__KG_") FROM SPACEXTABLE)
```

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

```
Out[22]: 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

- Explanation:
 - Listing the names of the booster versions which have carried the maximum payload mass.

2015 Launch Records

```
In [38]: %sql SELECT substr("Date",6,2) AS MONTH, "Landing_Outcome","Booster_Version","Launch_Site" FROM SPACEXTABLE WHERE "Landing_O
```



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

Out[38]:

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

- Explanation:
 - Listing the failed landing outcomes in drone ship, their booster versions and launch site names for the months in year 2015

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

```
In [40]: %sql SELECT "Landing_Outcome", COUNT(*) AS "Outcome_Count" FROM SPACEXTABLE GROUP BY "Landing_Outcome" HAVING "Date">'2010-06-04' AND "Date"<='2017-03-20'
```

* sqlite:///my_data1.db
Done.

Out[40]:

Landing_Outcome	Outcome_Count
No attempt	21
Success (drone ship)	14
Success (ground pad)	9
Failure (drone ship)	5
Controlled (ocean)	5
Uncontrolled (ocean)	2
Precluded (drone ship)	1

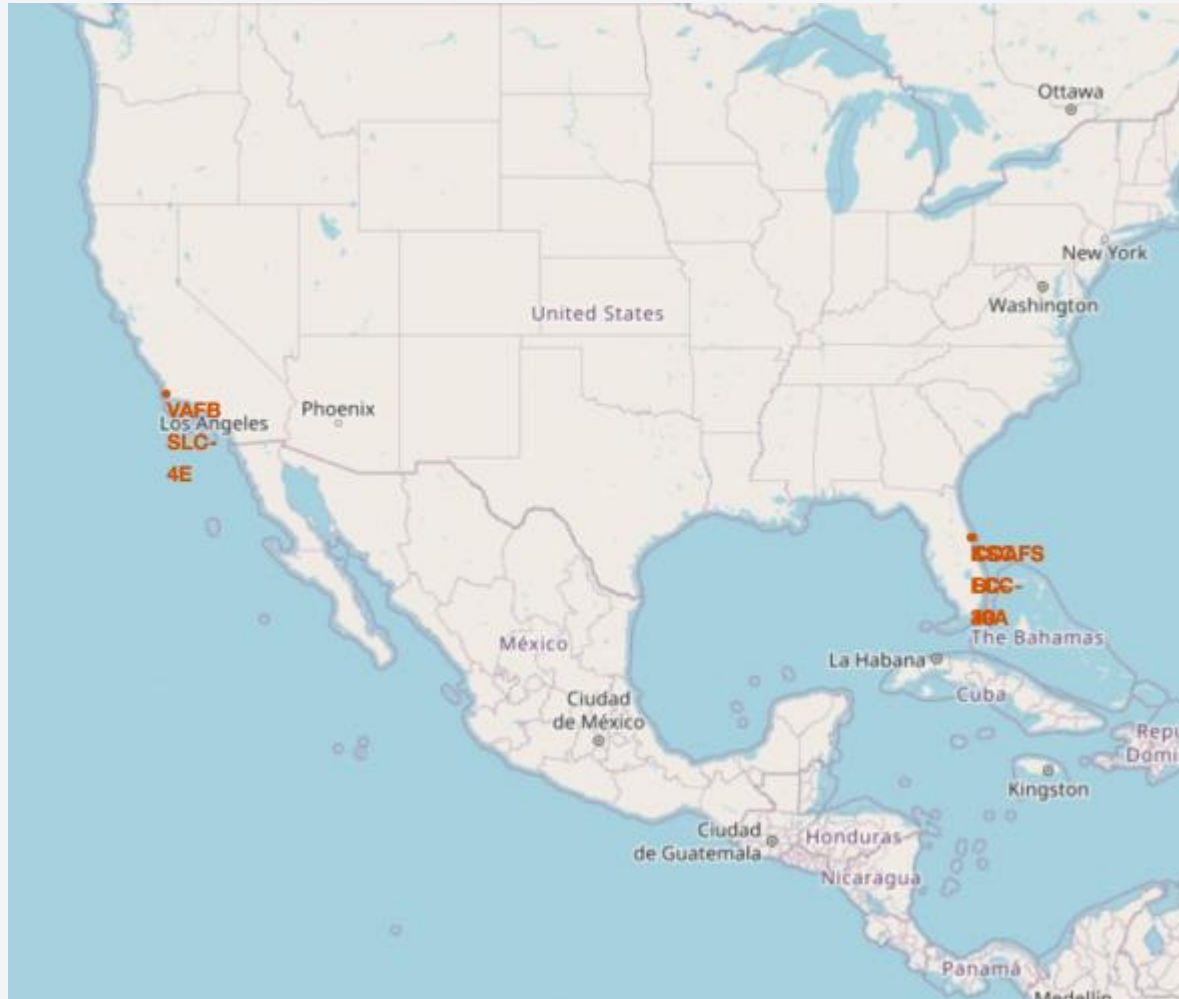
- Explanation: Ranking the count of landing outcomes (such as Failure (drone ship) or Success (ground pad)) between the date 2010-06-04 and 2017-03-20 in descending order

A satellite view of Earth from space, showing the curvature of the planet and city lights at night. The image is a composite of a solid blue rectangle on the left and a satellite photograph of Earth on the right. The Earth's surface is dark, with numerous bright yellow and orange lights representing cities and urban areas. The horizon of the Earth is visible, separating the dark surface from the deep blue of the atmosphere and the blackness of space.

Section 3

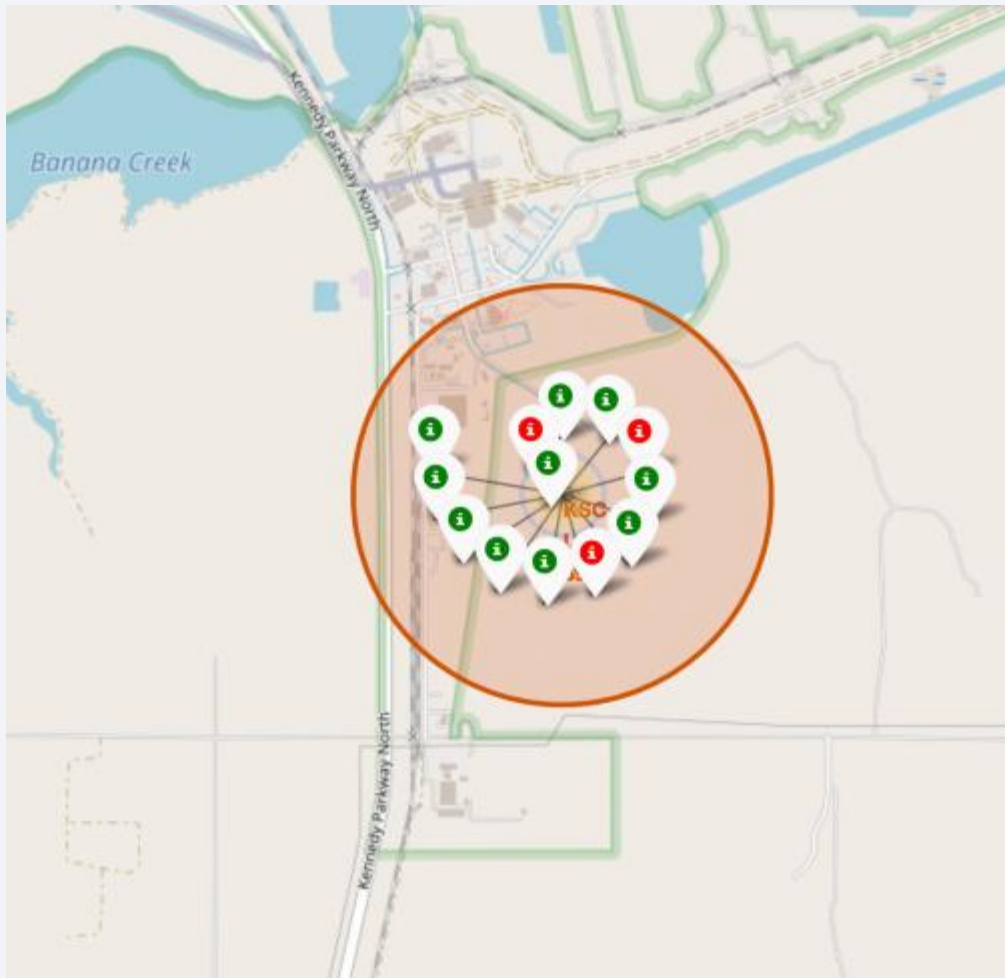
Launch Sites Proximities Analysis

<Folium Map Screenshot 1>



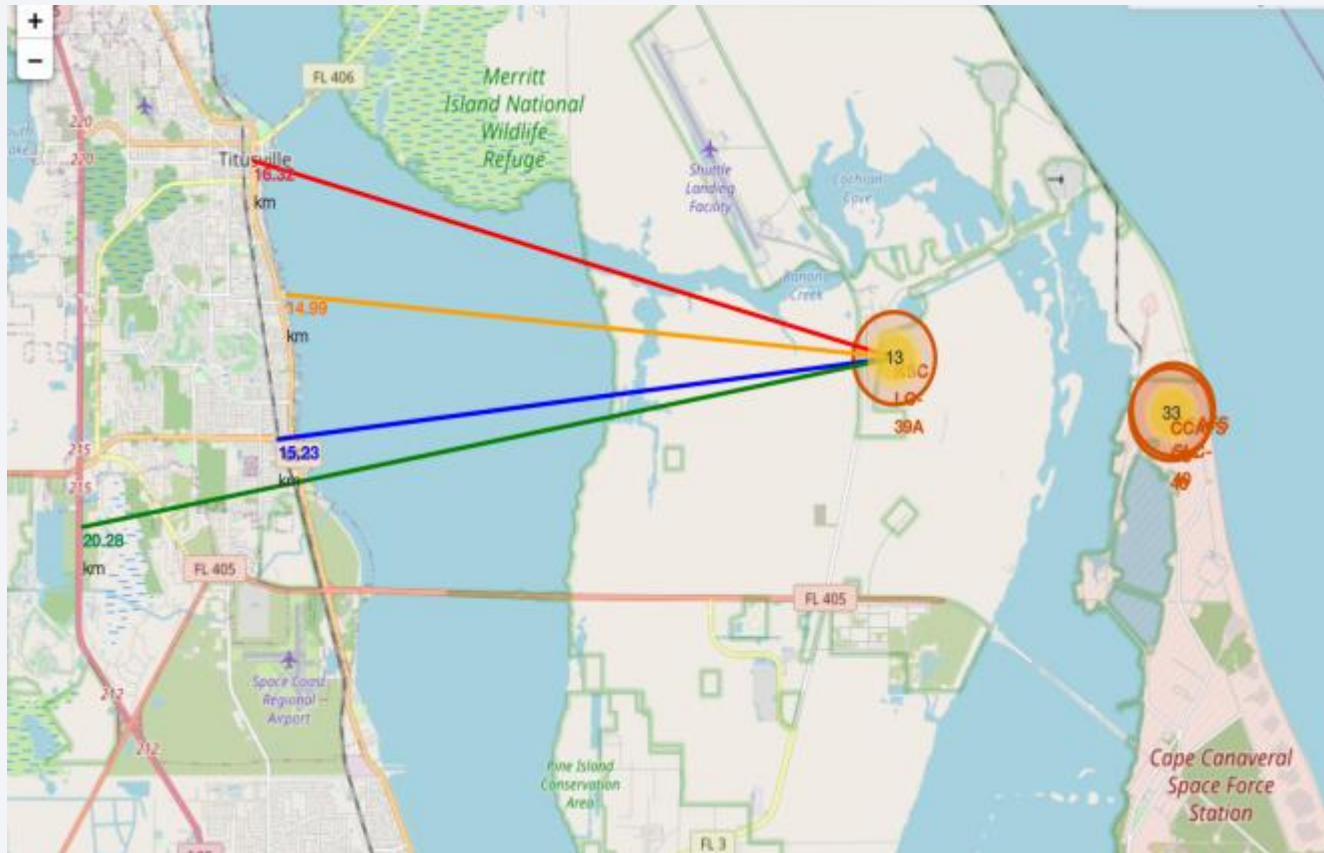
- Explanation:
 - Most of Launch sites are in proximity to the Equator line. The land is moving faster at the equator than any other place on the surface of the Earth. Anything on the surface of the Earth at the equator is already moving at 1670 km/hour. If a ship is launched from the equator it goes up into space, and it is also moving around the Earth at the same speed it was moving before launching. This is because of inertia. This speed will help the spacecraft keep up a good enough speed to stay in orbit.
 - All launch sites are in very close proximity to the coast, while launching rockets towards the ocean it minimises the risk of having any debris dropping or exploding near people.

<Folium Map Screenshot 2>



- Explanation:
 - From the colour-labeled markers we should be able to easily identify which launch sites have relatively high success rates.
 - Green Marker = Successful Launch
 - Red Marker = Failed Launch
 - Launch Site KSC LC-39A has a very high Success Rate.

<Folium Map Screenshot 3>



• Explanation:

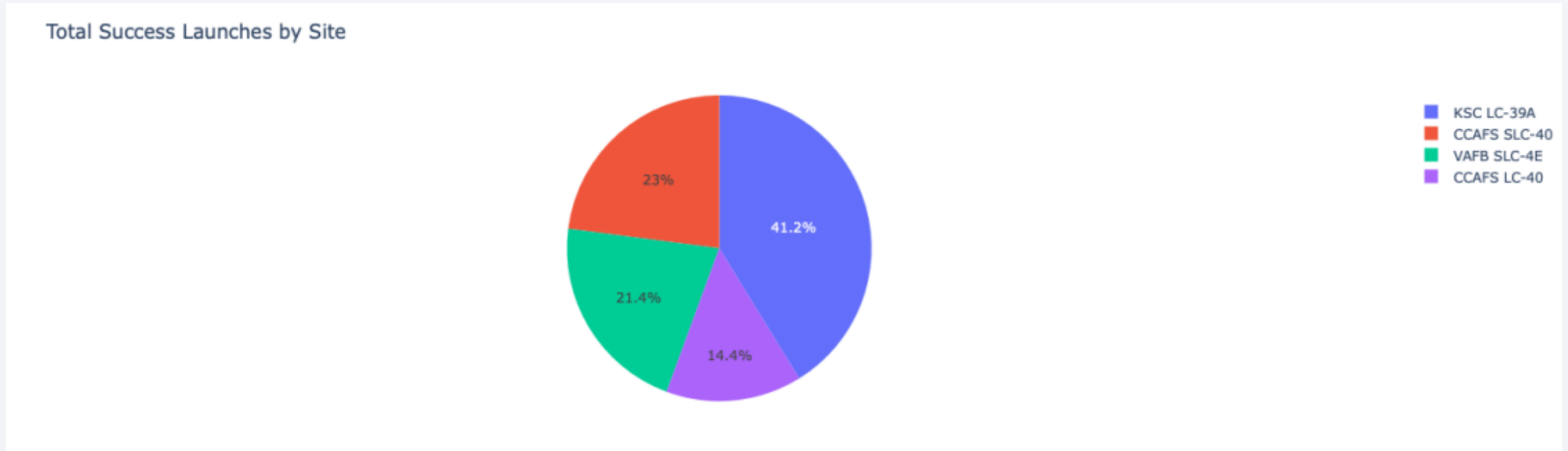
- From the visual analysis of the launch site KSC LC-39A we can clearly see that it is:
 - relative close to railway (15.23 km)
 - relative close to highway (20.28 km)
 - relative close to coastline (14.99 km)
- Also the launch site KSC LC-39A is relative close to its closest city Titusville (16.32 km).
- Failed rocket with its high speed can cover distances like 15-20 km in few seconds. It could be potentially dangerous to populated areas.



Section 4

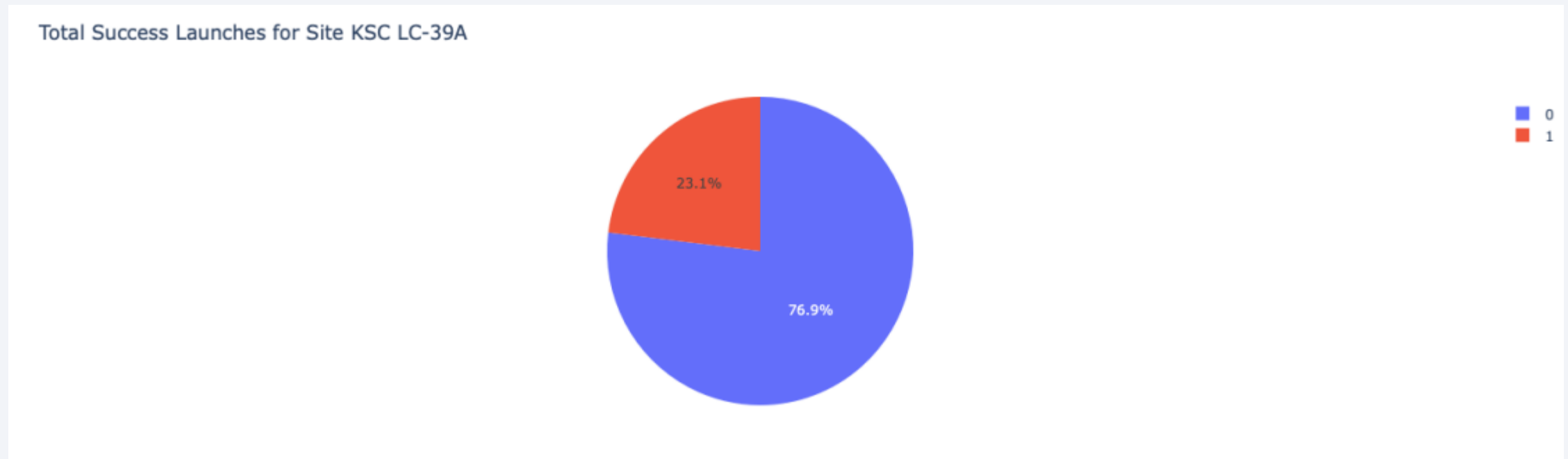
Build a Dashboard with Plotly Dash

<Dashboard Screenshot 1>



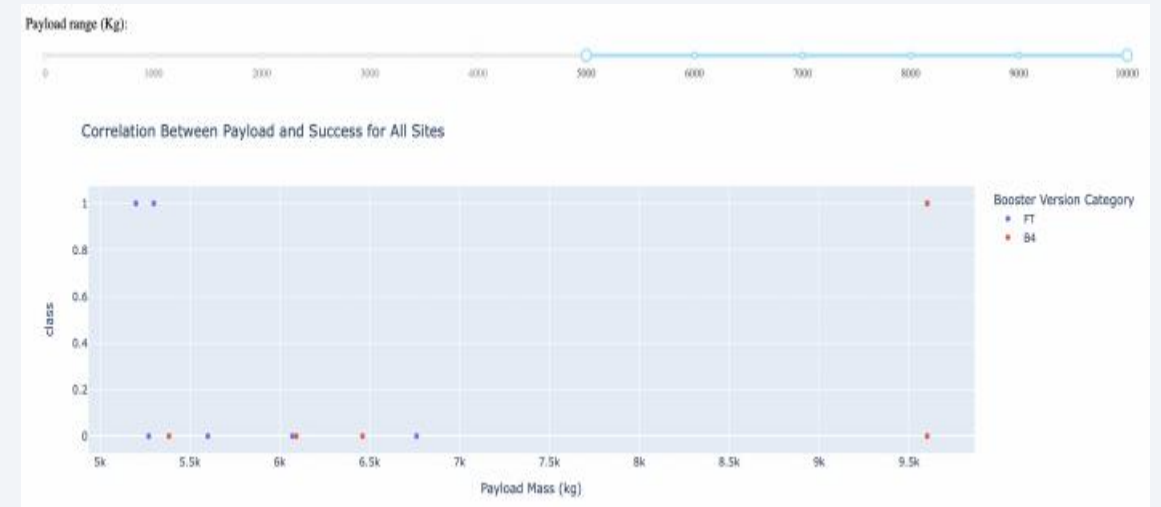
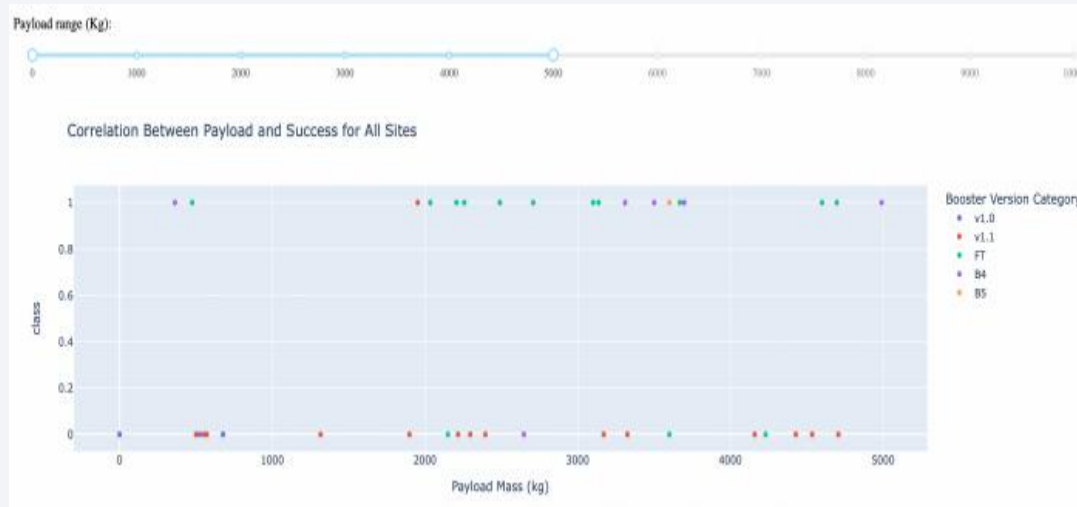
- Explanation:
 - The chart clearly shows that from all the sites, KSC LC-39A has the most successful launches.

<Dashboard Screenshot 2>



- Explanation:
 - KSC LC-39A has the highest launch success rate (76.9%) with 10 successful and only 3 failed landings.

<Dashboard Screenshot 3>



- Explanation:
 - The charts show that payloads between 2000 and 5500 kg have the highest success rate



Section 5

Predictive Analysis (Classification)

Classification Accuracy

- Explanation:
 - Based on the scores of the Test Set, we can not confirm which method performs best.
 - Same Test Set scores may be due to the small test sample size (18 samples). Therefore, we tested all methods based on the whole Dataset.
 - The scores of the whole Dataset confirm that the best model is the Decision Tree Model. This model has not only higher scores, but also the highest accuracy.

	LogReg	SVM	Tree	KNN
Jaccard_Score	0.800000	0.800000	0.800000	0.800000
F1_Score	0.888889	0.888889	0.888889	0.888889
Accuracy	0.833333	0.833333	0.666667	0.833333

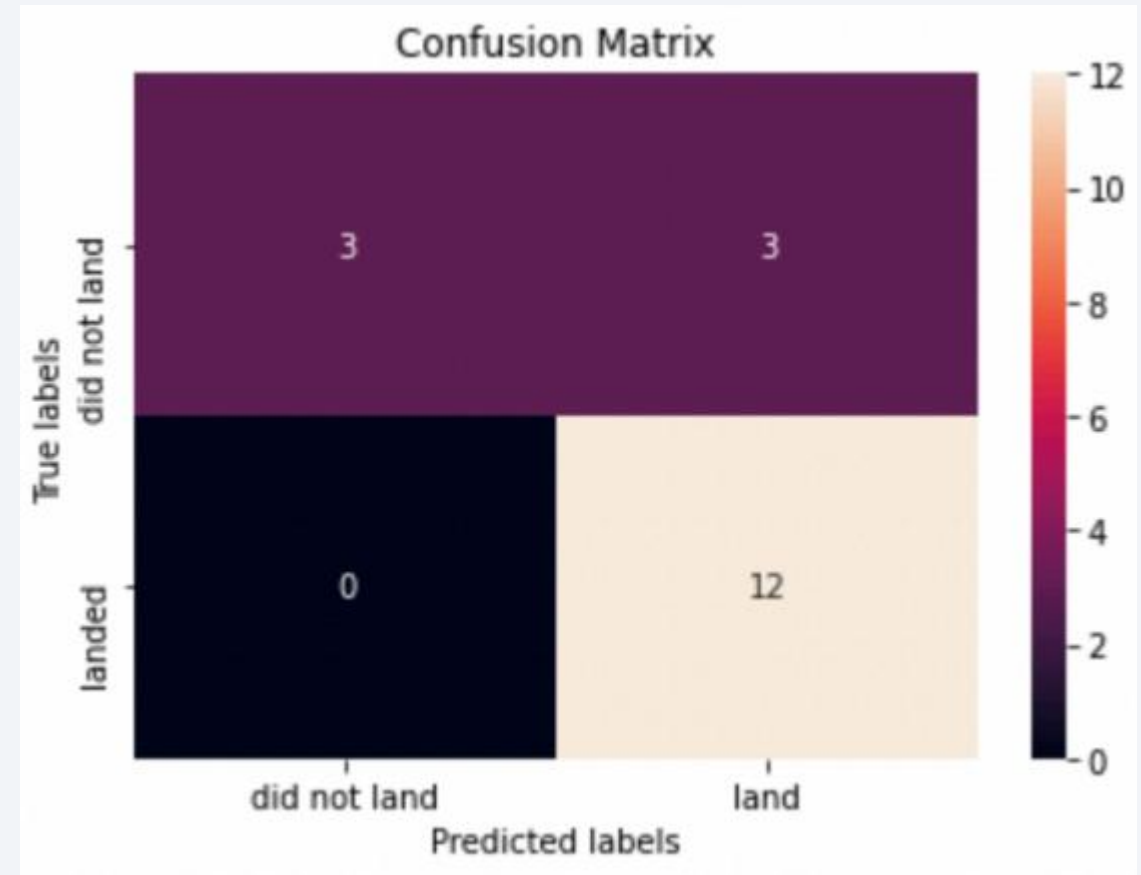
Scores and Accuracy of the Test Set

	LogReg	SVM	Tree	KNN
Jaccard_Score	0.833333	0.845070	0.882353	0.819444
F1_Score	0.909091	0.916031	0.937500	0.900763
Accuracy	0.866667	0.877778	0.911111	0.855556

Scores and Accuracy of the Entire Data Set

Confusion Matrix

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



Conclusions

- Decision Tree Model is the best algorithm for this dataset.
- Launches with a low payload mass show better results than launches with a larger payload mass.
- Most of launch sites are in proximity to the Equator line and all the sites are in very close proximity to the coast.
- The success rate of launches increases over the years.
- KSC LC-39A has the highest success rate of the launches from all the sites.
- Orbits ES-L1, GEO, HEO and SSO have 100% success rate

Appendix

- Special Thanks to: Instructors, Coursera and IBM

Thank you!

