



IBM Developer  
SKILLS NETWORK

# Winning Space Race with Data Science

Nguyen Quoc Hung  
17/7/2024



# Outline

---

- Executive Summary
- Introduction
- Methodology
- Results
- Conclusion

# Executive Summary

---

- Summary of methodologies
  - Data Collection through API and Web Scraping
  - Data Wrangling
  - Exploratory Data Analysis with SQL and Data Visualization
  - Interactive Visual Analytics With Folium
  - Machine Learning Prediction
- Summary of all results
  - Exploratory Data Analysis result
  - Interactive analytics in screenshots
  - Predictive Analytics result

# Introduction

---

- In this capstone, we will predict if the success of the Falcon 9 rocket's first stage landing. SpaceX offers Falcon 9 launches at a significantly lower cost of 62 million dollars, compared to other providers charging upwards of 165 million dollars. This cost efficiency is largely due to SpaceX's ability to reuse the first stage of their rocket
- Predicting whether the first stage will land successfully helps estimate launch costs, providing valuable insight for companies looking to compete with SpaceX. This information can help these companies strategize their bids more effectively



Section 1

# Methodology

# Methodology

---

## Executive Summary

- Data collection methodology:
  - Data was collected using SpaceX API and web scraping from wikipedia
- Perform data wrangling
  - One-hot encoding was applied to categorical features
- Perform exploratory data analysis (EDA) using visualization and SQL
- Perform interactive visual analytics using Folium and Plotly Dash
- Perform predictive analysis using classification models
  - Standardize the data
  - Find the method performs best

# Data Collection – SpaceX API

- The API used is [api.spacexdata.com/v4/rockets/](https://api.spacexdata.com/v4/rockets/)
- The API provide 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 all 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
--	--	--	--	--	--	--	--	--	--	--	--	--		--	--	--	--	--
89	86	2020-09-03	Falcon 9	15600.0	VLEO	KSC LC 39A	True ASDS	2	True	True	True	5e9e3032383ecb6bb234e7ca	5.0	12	B1060	-80.603956	28.608058	
90	87	2020-10-06	Falcon 9	15600.0	VLEO	KSC LC 39A	True ASDS	3	True	True	True	5e9e3032383ecb6bb234e7ca	5.0	13	B1058	-80.603956	28.608058	
91	88	2020-10-18	Falcon 9	15600.0	VLEO	KSC LC 39A	True ASDS	6	True	True	True	5e9e3032383ecb6bb234e7ca	5.0	12	B1051	-80.603956	28.608058	
92	89	2020-10-24	Falcon 9	15600.0	VLEO	CCSFS SLC 40	True ASDS	3	True	True	True	5e9e3033383ecbb9e534e7cc	5.0	12	B1060	-80.577366	28.561857	
93	90	2020-11-05	Falcon 9	3681.0	MEO	CCSFS SLC 40	True ASDS	1	True	False	True	5e9e3032383ecb6bb234e7ca	5.0	8	B1062	-80.577366	28.561857	

# Data Collection - Scraping

- The data is scraped from [en.wikipedia.org/w/index.php?title=List of Falcon 9 and Falcon Heavy launches&oldid=1027686922](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 all rows of the data:

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



# Data Wrangling

---

- We create a landing outcome label from Outcome column and then assign it to class column
- We end up with 91 rows or instances and 18 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	Class
0	1	2010-06-04	Falcon 9	6123.547647	LEO	CCSFS SLC 40	None None	1	False	False	False	NaN	1.0	0	B0003	-80.577366	28.561857	0
1	2	2012-05-22	Falcon 9	525.000000	LEO	CCSFS SLC 40	None None	1	False	False	False	NaN	1.0	0	B0005	-80.577366	28.561857	0
2	3	2013-03-01	Falcon 9	677.000000	ISS	CCSFS SLC 40	None None	1	False	False	False	NaN	1.0	0	B0007	-80.577366	28.561857	0
3	4	2013-09-29	Falcon 9	500.000000	PO	VAFB SLC 4E	False Ocean	1	False	False	False	NaN	1.0	0	B1003	-120.610829	34.632093	0
4	5	2013-12-03	Falcon 9	3170.000000	GTO	CCSFS SLC 40	None None	1	False	False	False	NaN	1.0	0	B1004	-80.577366	28.561857	0

# EDA with Data Visualization

---

- Pandas and NumPy use to derive basic information about the data, which include:
  - The number of launches on each launch site
  - The number of occurrence of each orbit
  - The number and occurrence of each mission outcome
- SQL use 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

# Build an Interactive Map with Folium

---

- 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

# Build a Dashboard with Plotly Dash

---

- Build an interactive dashboard that contains pie charts and scatter plots to analyze data with the Plotly Dash Python library.
- Calculate distances on an interactive map by writing Python code using the Folium library.
- Generate interactive maps, plot coordinates, and mark clusters by writing Python code using the Folium library.
- Build a dashboard to analyze launch records interactively with Plotly Dash.
- Build an interactive map to analyze the launch site proximity with Folium.

# Predictive Analysis (Classification)

---

- Split the data into training testing data
- Train different classification models
- Optimize the Hyperparameter grid search
- Evaluate the models based on their accuracy scores and confusion matrix



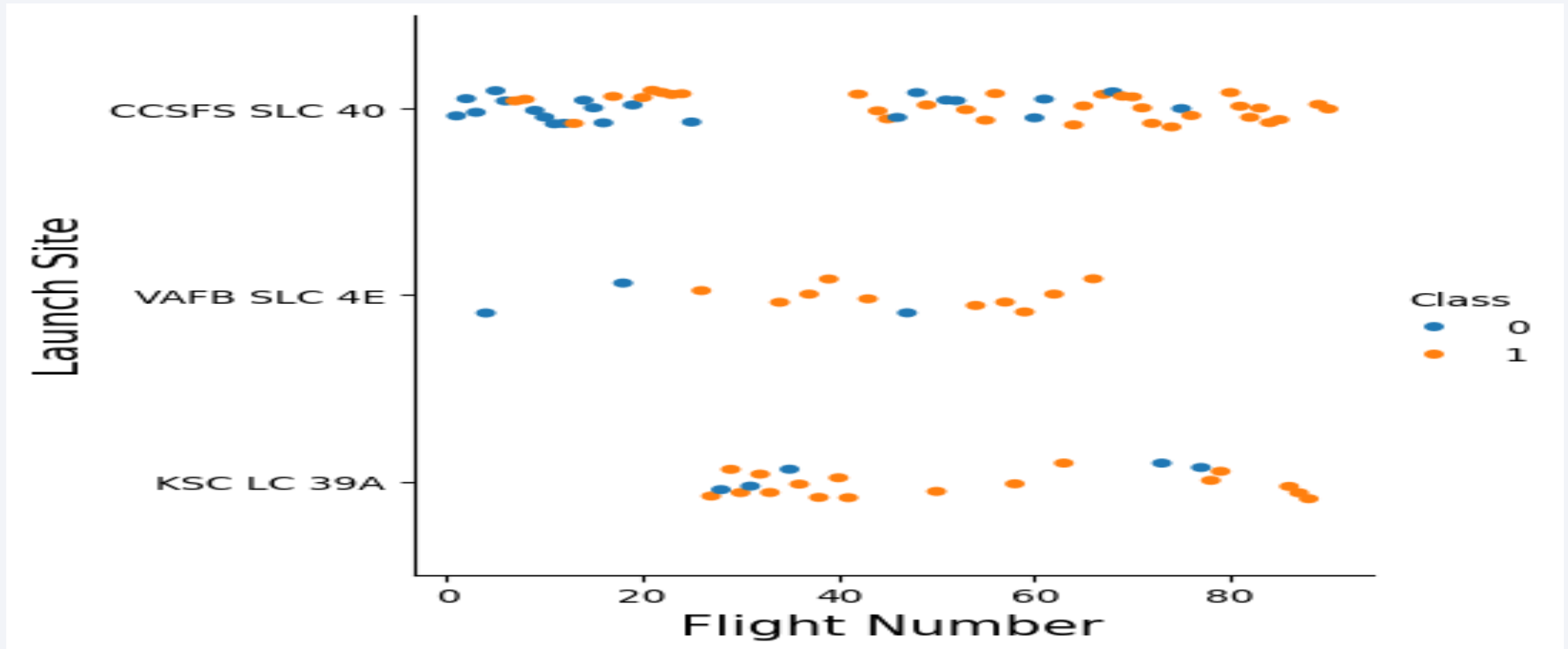
The background of the slide is an abstract composition. It features a dark blue base color. Overlaid on this are numerous diagonal streaks in shades of red and cyan. A faint, light blue grid pattern is also visible, particularly in the lower half of the image. The overall effect is dynamic and technological.

Section 2

# Insights drawn from EDA

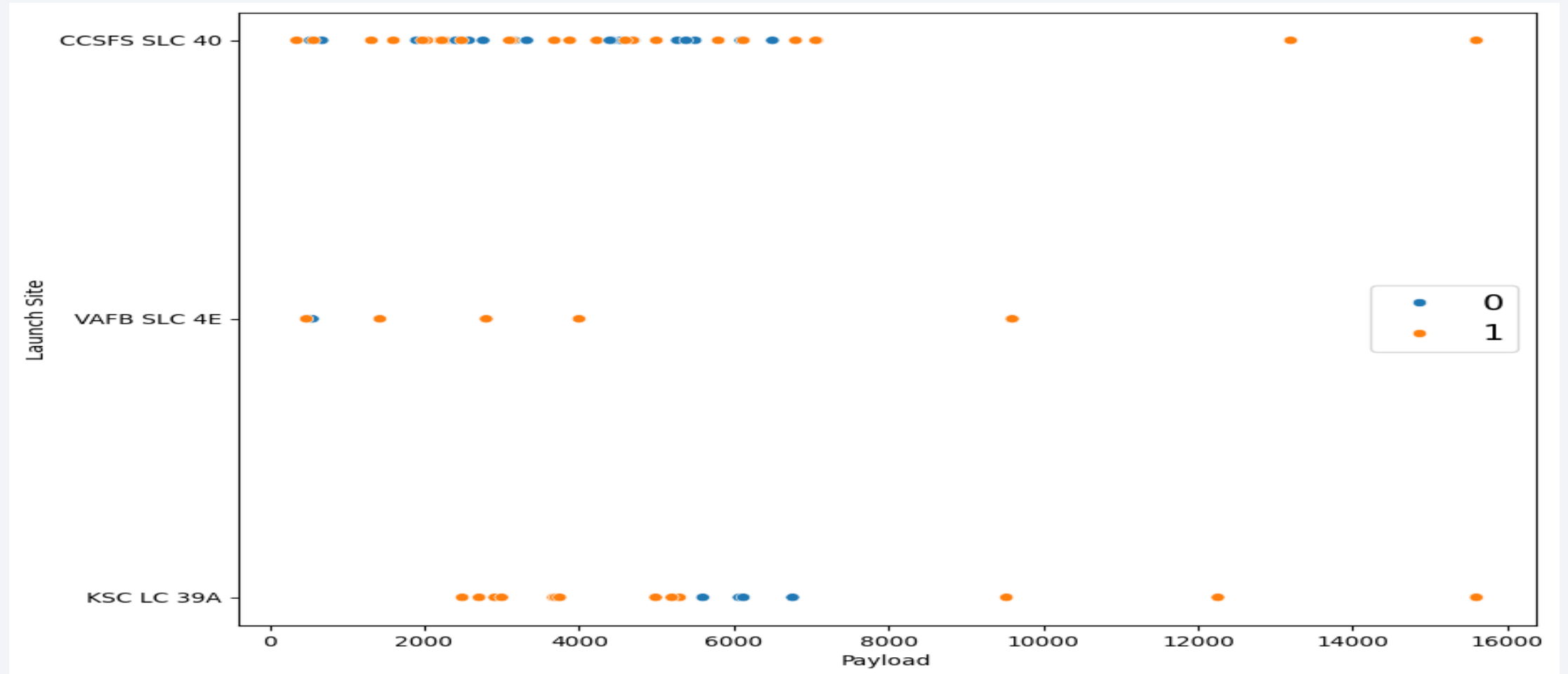


# Flight Number vs. Launch Site



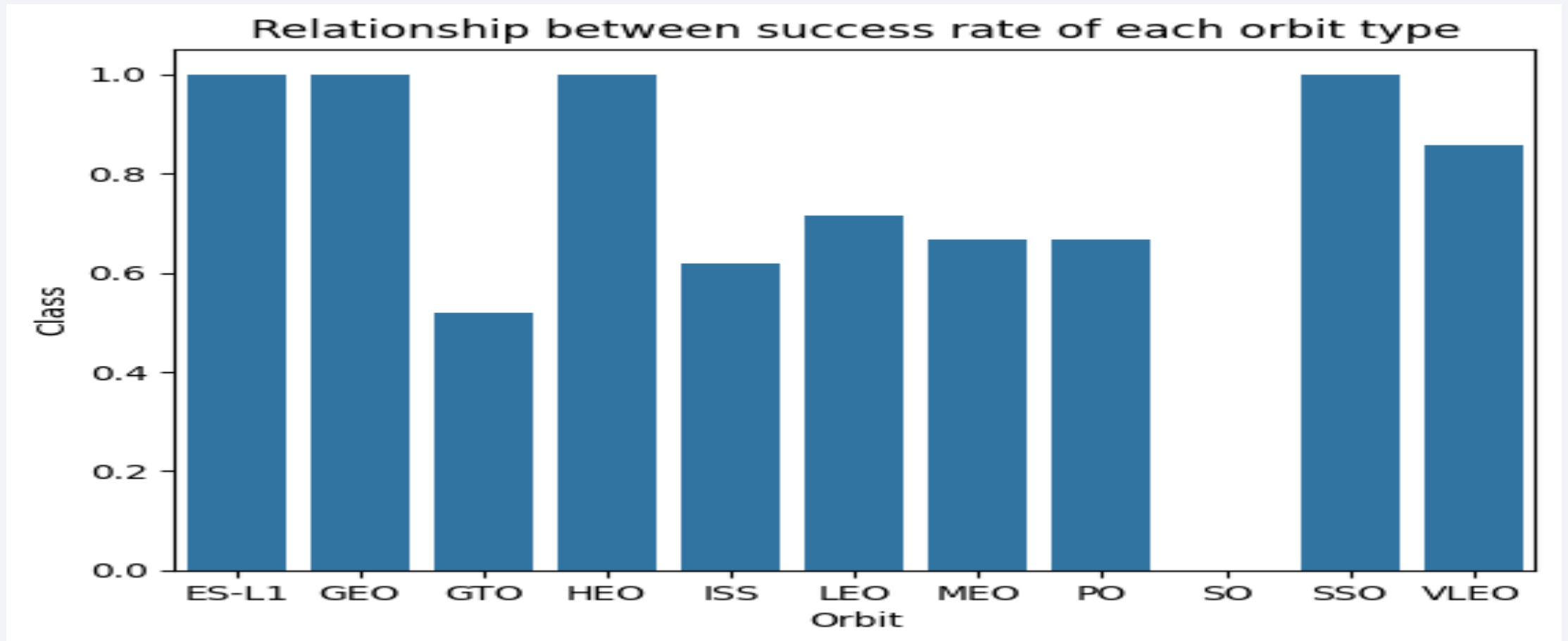
scatter plot of Flight Number vs. Launch Site

# Payload vs. Launch Site



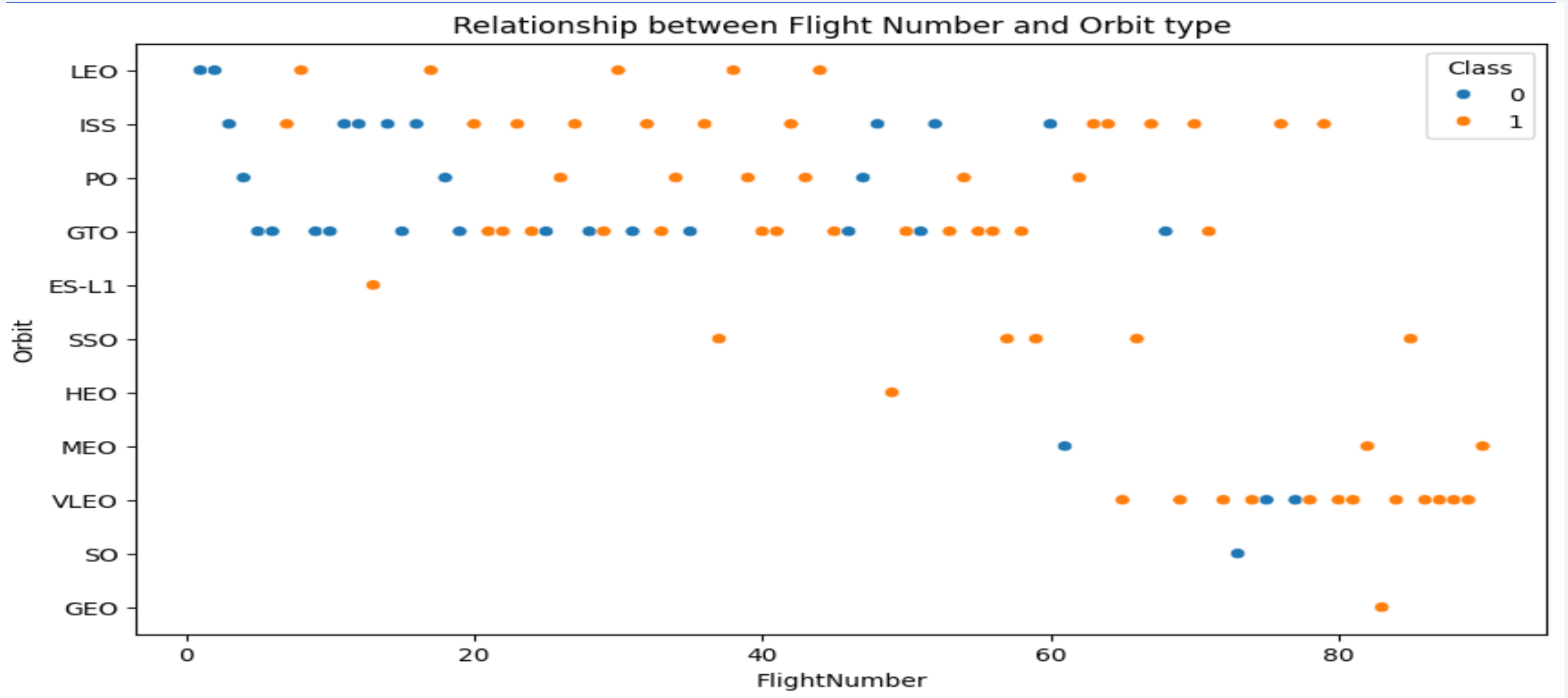
scatter plot of Payload vs. Launch Site

# Success Rate vs. Orbit Type



bar chart for the success rate of each orbit type

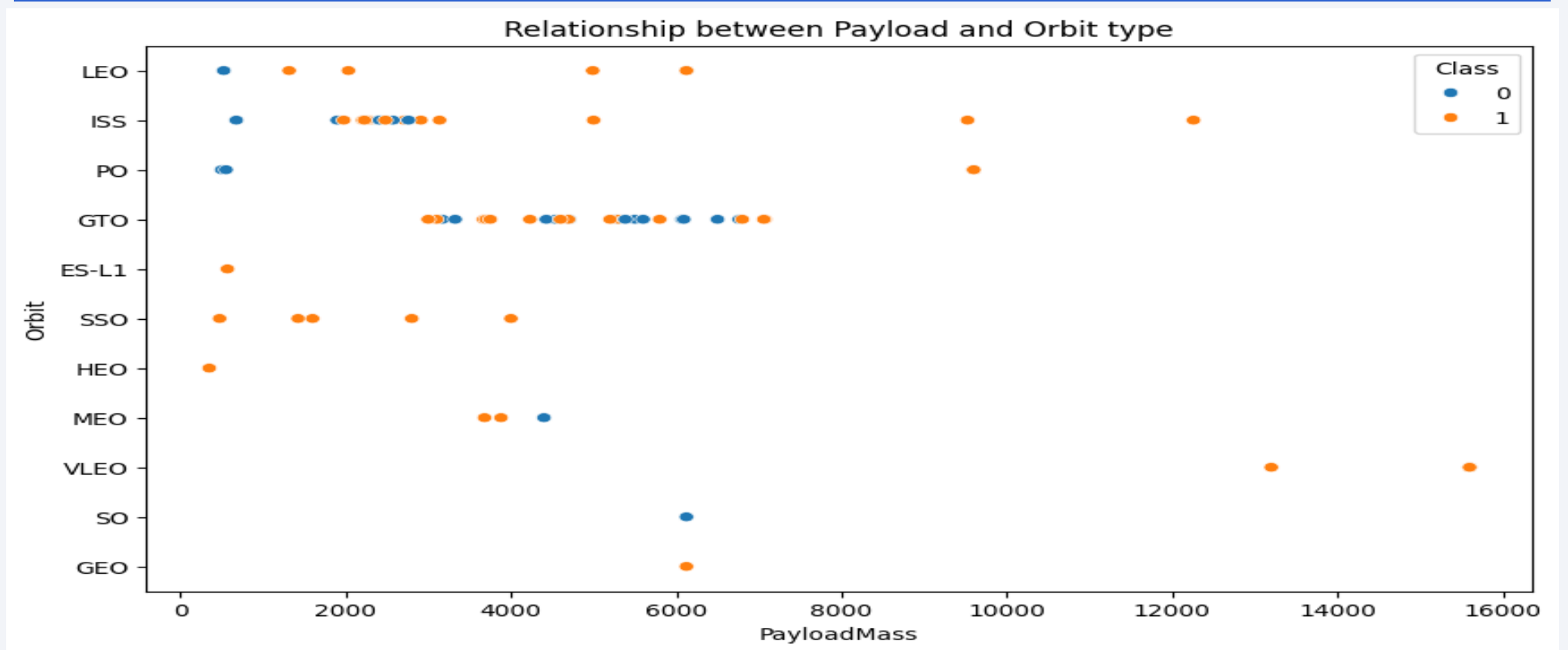
# Flight Number vs. Orbit Type



scatter point of Flight number vs. Orbit type

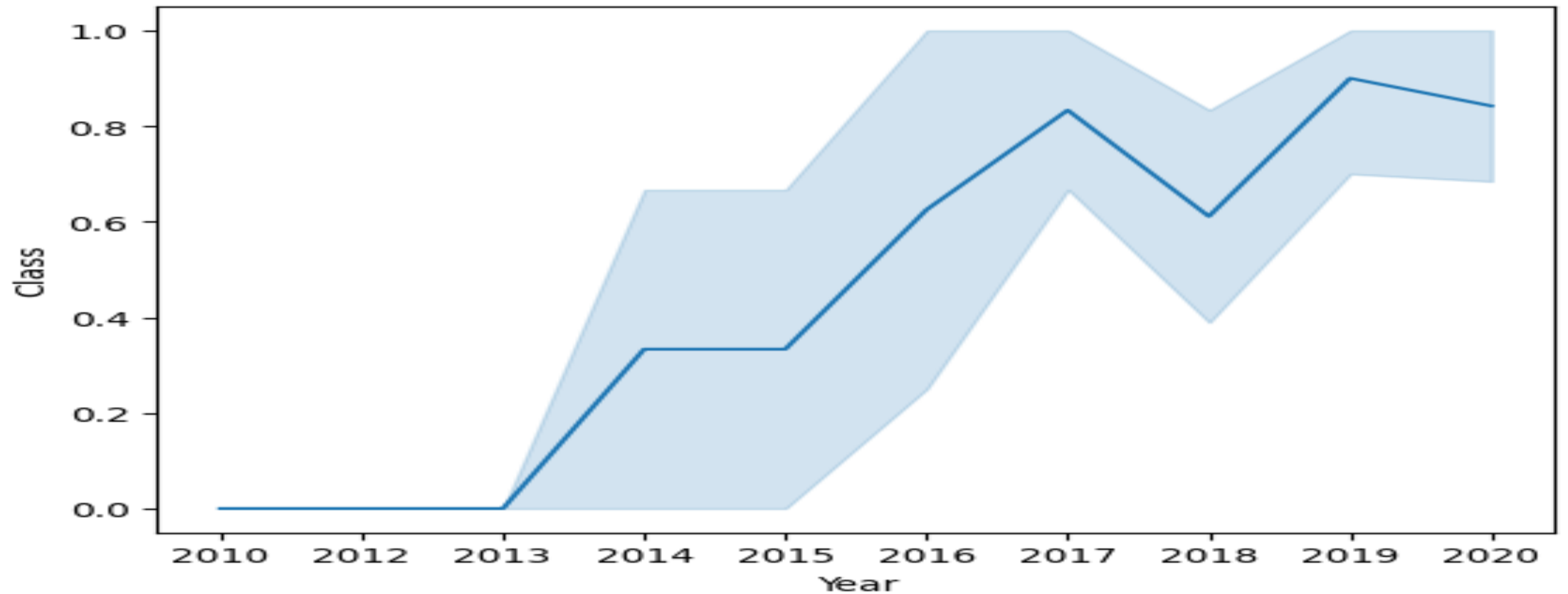


# Payload vs. Orbit Type



scatter point of payload vs. orbit type

# Launch Success Yearly Trend



line chart of yearly average success rate

# All Launch Site Names

---

Names of the unique launch sites

```
%sql select distinct(Launch_Site)from SPACEXTBL
```

```
* sqlite:///my\_data1.db
```

```
Done.
```

```
Launch_Site
```

```
CCAFS LC-40
```

```
VAFB SLC-4E
```

```
KSC LC-39A
```

```
CCAFS SLC-40
```

# Launch Site Names Begin with 'CCA'

5 records where launch sites begin with `CCA`

```
%sql select * from SPACEXTBL where Launch_Site like 'CCA%' limit 5
```

```
* sqlite:///my\_data1.db
```

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	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

# Total Payload Mass

---

Total payload carried by boosters from NASA

```
%sql select sum(PAYLOAD_MASS_KG_) from SPACEXTBL where Customer='NASA (CRS)'
```

```
* sqlite:///my\_data1.db
```

```
Done.
```

```
sum(PAYLOAD_MASS_KG_)
```

```
45596
```



# Average Payload Mass by F9 v1.1

---

Average payload mass carried by booster version F9 v1.1

```
%sql select avg(PAYLOAD_MASS_KG_) from SPACEXTBL where Booster_Version like '%F9 v1.1%'
```

```
* sqlite:///my\_data1.db
```

```
Done.
```

```
avg(PAYLOAD_MASS_KG_)
```

```
2534.6666666666665
```

# First Successful Ground Landing Date

---

Dates of the first successful landing outcome on ground pad

```
%sql select min(Date) from SPACEXTBL where Mission_Outcome = 'Success'
```

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

```
Done.
```

```
min(Date)
```

```
2010-06-04
```

# 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 distinct(Booster_Version) from SPACEXTBL where PAYLOAD_MASS_KG_ between 4000 and 6000

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

Booster_Version
F9 v1.1
F9 v1.1 B1011
F9 v1.1 B1014
F9 v1.1 B1016
F9 FT B1020
F9 FT B1022
F9 FT B1026
F9 FT B1030
F9 FT B1021.2
F9 FT B1032.1
F9 B4 B1040.1
F9 FT B1031.2
F9 B4 B1043.1
F9 FT B1032.2
F9 B4 B1040.2
F9 B5 B1046.2
F9 B5 B1047.2
F9 B5 B1046.3
F9 B5B1054
F9 B5 B1048.3
F9 B5 B1051.2
F9 B5B1060.1
F9 B5 B1058.2
F9 B5B1062.1

# Total Number of Successful and Failure Mission Outcomes

---

Total number of successful and failure mission outcomes

```
%sql select count(*), Mission_Outcome from SPACEXTBL group by Mission_Outcome
```

```
* sqlite:///my\_data1.db
```

```
Done.
```

count(*)	Mission_Outcome
1	Failure (in flight)
98	Success
1	Success
1	Success (payload status unclear)

# 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)
```

```
* sqlite:///my_data1.db  
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 strftime('%m', Date) AS Month, Landing_Outcome, Booster_Version, Launch_Site FROM SPACEXTBL WHERE Landing_Outcome = 'Failure (drone ship)' AND strftime('%Y', Date) = '2015';
```

```
* sqlite:///my\_data1.db
```

```
Done.
```

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

# 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 descending order

```
%sql select count(*) as cnt, Landing_Outcome from SPACEXTBL group by Landing_Outcome order by cnt DESC
```

```
* sqlite:///my\_data1.db
```

```
Done.
```

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

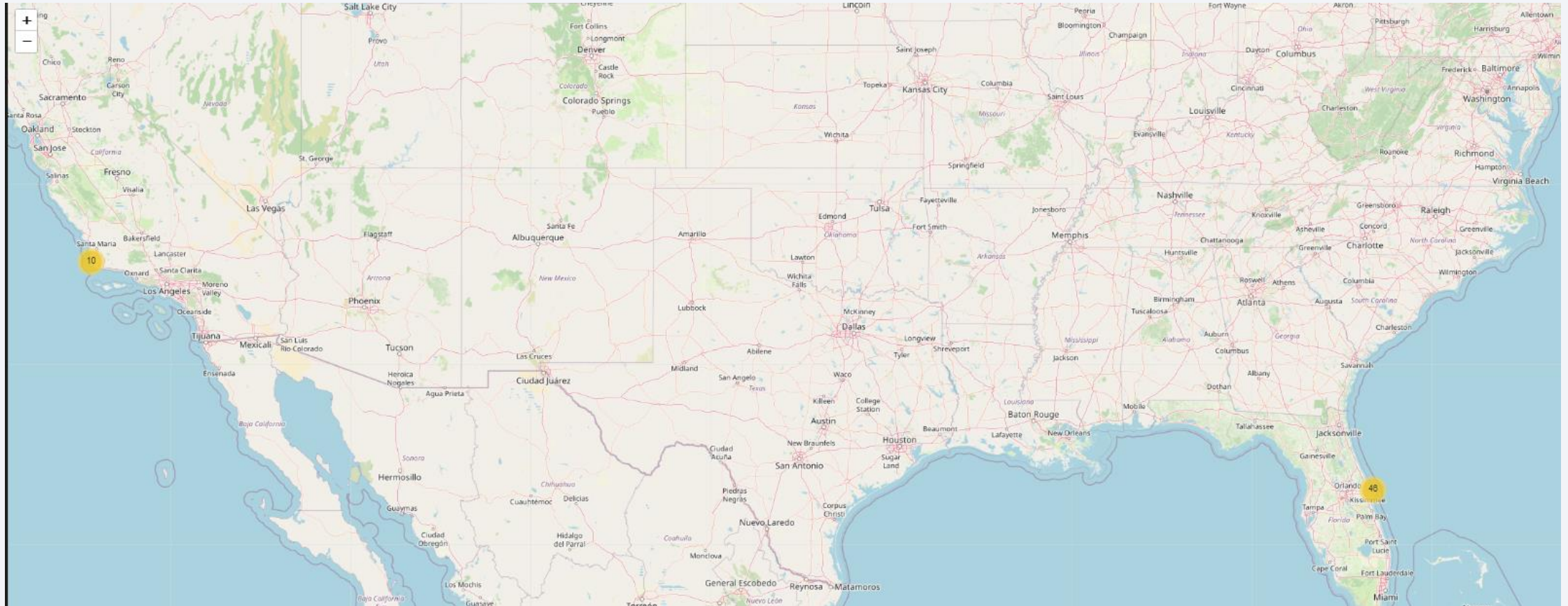
A satellite view of Earth from space, showing the curvature of the planet and city lights at night. The background is a deep blue gradient.

Section 3

# Launch Sites Proximities Analysis

# Folium Map

All launch site on Map

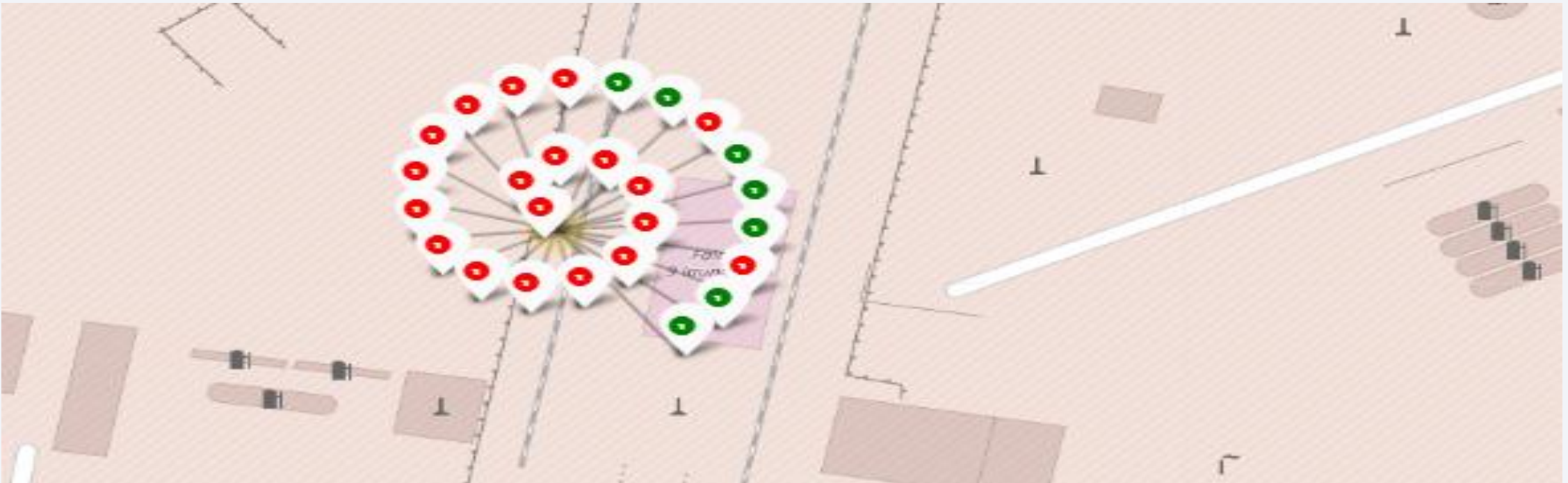




# Folium Map

---

- The succeeded launches and failed launches for each site on map. Each green tag represents a successful launch while each red tag represents a failed launch



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.

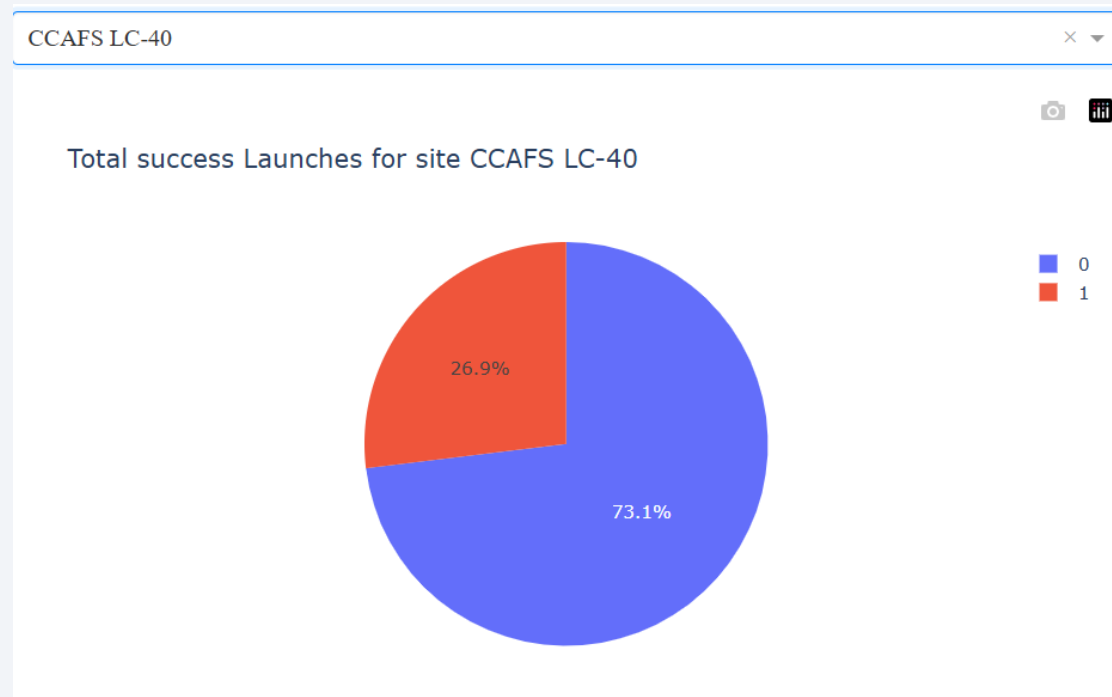
Section 4

# Build a Dashboard with Plotly Dash

# Dashboard

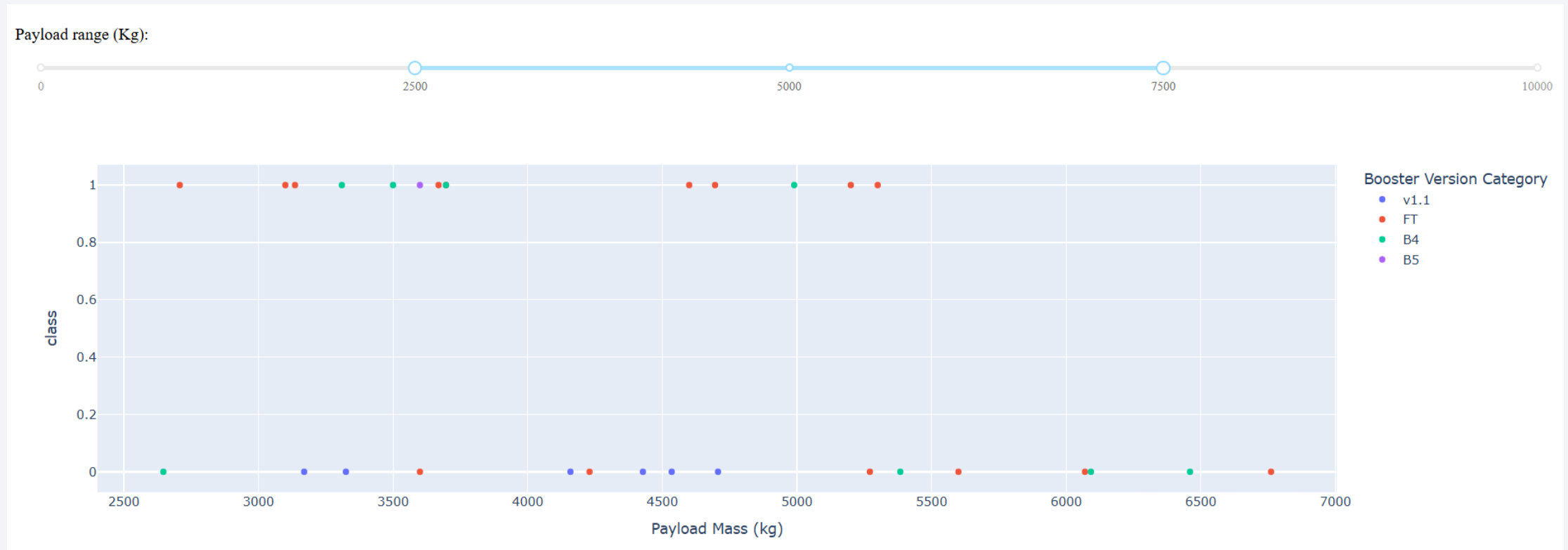
---

- The picture below shows a pie chart when launch site CCAFS LC-40 is chosen.
- 0 represents failed launches while 1 represents successful launches. We can see
- that 73.1% of launches done at CCAFS LC-40 are failed launches



# Dashboard

- The picture below shows a scatterplot when the payload mass range is set to be from 2500kg to 7500kg.
- Class 0 represents failed launches while class 1 represents successful launches





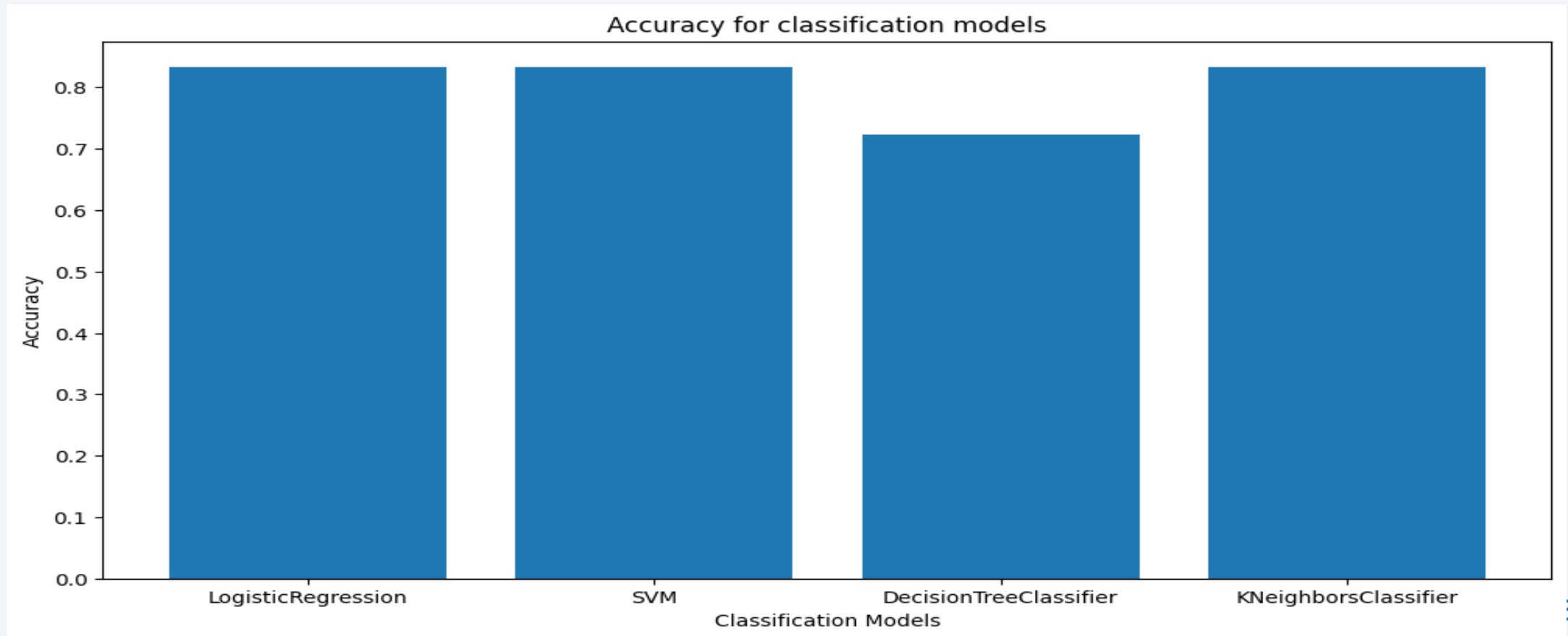
Section 5

# Predictive Analysis (Classification)

# Classification Accuracy

---

- Putting the results of all 4 models side by side, we can see that they all share the same accuracy score when test on test set.



# Conclusions

---

- • In this project, we try to predict if the first stage of a given Falcon 9 launch will
- land in order to determine the cost of a launch.
- • Each feature of a Falcon 9 launch, such as its payload mass or orbit type, may
- affect the mission outcome in a certain way.
- • Several machine learning algorithms are employed to learn the patterns of past
- Falcon 9 launch data to produce predictive models that can be used to predict the
- outcome of a Falcon 9 launch.
- • The predictive model produced by decision tree algorithm performed the best
- among the 4 machine learning algorithms employed

Thank you!

