

Winning Space Race with Data Science

By, Advait Krishna



Outline

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

Executive Summary

Summary of methodologies

The project involves a comprehensive data-driven journey, starting with the collection of data through API interactions, which allows for the retrieval of relevant information from various sources. Subsequently, web scraping is employed as another robust method for data collection, expanding the scope of available data. Following the data gathering phase, the focus shifts to data wrangling, where the collected raw data is organized and cleaned to ensure its suitability for analysis. The exploration begins with Exploratory Data Analysis (EDA) using SQL queries to extract meaningful insights from structured datasets. This is complemented by visual exploration techniques, utilizing data visualization tools to present patterns and trends effectively. Further enhancing the analytical depth, interactive visual analytics with Folium are employed, providing an immersive experience in geospatial data representation. The project then advances into the realm of machine learning prediction, leveraging the insights gained from previous stages to develop models for forecasting and decision-making. This comprehensive approach ensures a systematic and thorough exploration of the dataset, leading to valuable insights and informed outcomes.

Summary of all results

The success rates of launch sites exhibit a positive correlation with increasing flight numbers, indicating a trend where a higher flight number corresponds to a greater likelihood of success. Among specific orbital categories, including ES-L1, GEO, HEO, SSO, and VLEO, the success rates are notably high, suggesting the efficiency and reliability of launches into these orbits. Notably, missions with heavy payloads, particularly those destined for Polar, LEO, and ISS, demonstrate higher rates of successful landings. Analyzing the temporal aspect, the success rate has shown a consistent upward trajectory from 2013 to 2020, reflecting an overall improvement in launch outcomes. Geographically, all launch sites are strategically located near the equator and in close proximity to the American coast, specifically in Florida and California. Despite their coastal positioning, launch sites maintain a considerable distance from highways, railways, and urban areas. Furthermore, the success rates vary based on payload weight, with low-weighted payloads (0-4000kg) exhibiting a higher success rate compared to heavier payloads (4000-10000kg). Among the launch sites, KSC LC – 39A stands out with a remarkable success rate of 76.9% and a failure rate of 23.1%, marking it as the site with the highest success rate among the three considered launch locations.

3

Introduction

Project background and context

SpaceX promotes Falcon 9 rocket launches on its website at a price of 62 million dollars, significantly lower than the costs associated with other providers, which can reach upwards of 165 million dollars per launch. A substantial portion of this cost advantage stems from SpaceX's innovative ability to reuse the first stage of the rocket. Consequently, the success or failure of landing the first stage plays a crucial role in determining the overall launch cost. This information holds strategic significance for potential competitors seeking to bid against SpaceX for rocket launches. The primary objective of this project is to develop a machine learning pipeline capable of predicting the successful landing of the first stage, thereby providingvaluable insights into the cost estimation for rocket launches.

Problems you want to find answers

- What influences the successful landing of a rocket?
- The interplay of different factors that contribute to the success rate of a rocket's landing.
- What operational conditions must be established to guarantee the success of a landing program?



Methodology

Executive Summary

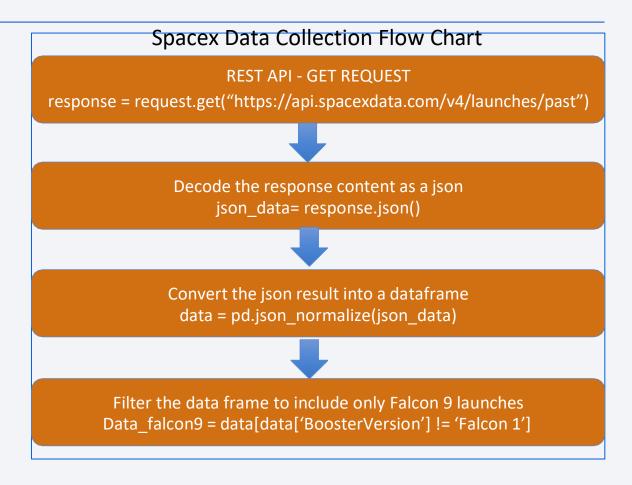
- Data was collection through the SpaceX REST API and Web Scrapping from Wikipedia website
- Perform data wrangling
- Perform exploratory data analysis (EDA) using visualization and SQL
- Perform interactive visual analytics using Folium and Plotly Dash
- Perform predictive analysis using classification models

Data Collection

• Data Collection Process and Flowcharts are described in following slides

Data Collection - SpaceX API

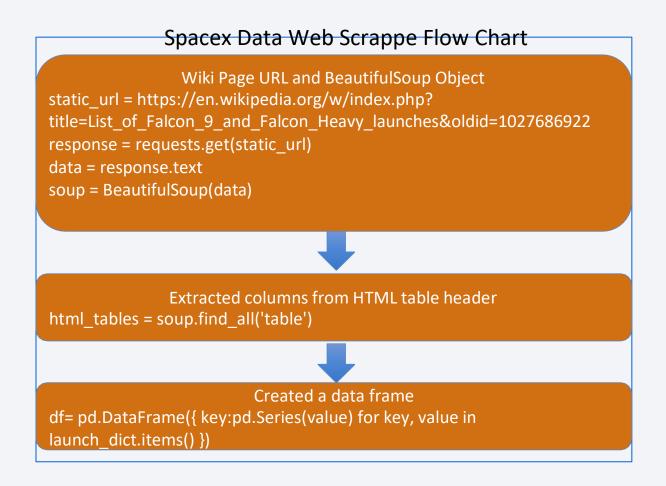
 First, we collected data from SpaceX API with get request method, then we clean the requested data as we needed and then perform some data wrangling and formatting on it.



The GitHub URL of the SpaceX API calls notebook

Data Collection - Scraping

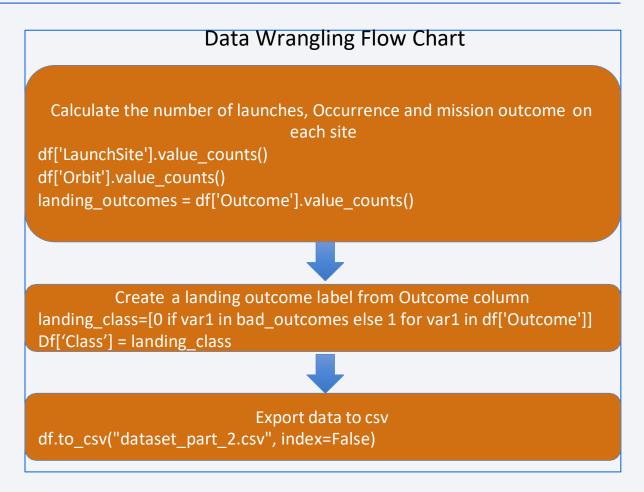
- First, we request the Falcon9 launch Wiki page from its URL and create a BeautifulSoup object from the HTML response.
- Then extracted all column names from the HTML table header.
- Finally, created a data frame by parsing the launch HTML Table.



GitHub URL link of the web scraping notebook:

Data Wrangling

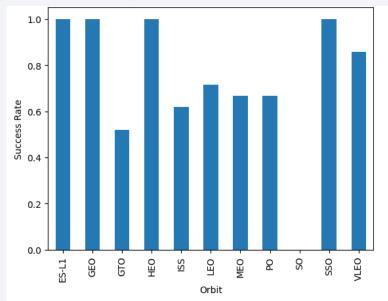
- To find data pattern and training label, we perform Exploratory Data Analysis(EDA).
- We calculated each site's number of launch, number and occurrence and mission outcome of each orbit.
- Finally, created a landing outcome label from Outcome columns and exported the csv data.



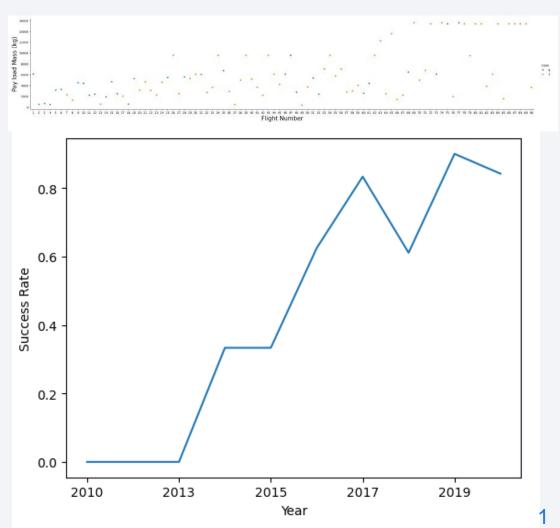
GitHub URL of data wrangling related notebooks:

EDA with Data Visualization

 The data exploration phase involved visualizing relationship between Flight Number and Launch sites, Payload Mass and Launch Site, success rate of each orbit type, Flight Number and Orbit type and Yearly Launch Success.



 GitHub URL of EDA with data visualization notebook:



https://github.com/sureshmanjoshi/testrepo/blob/main/jupyter-labs-eda-dataviz.ipynb

EDA with SQL

- We used SQlite database to load Spacex data
- We used SQL queries to get following information from the data:
 - Finding the names of unique launch sites in the space mission
 - Find out the total payload mass carried by boosters launched by NASA
 - Find out average payload mass carried by booster version F9 v1.1
 - Find the first successful landing outcome in ground pad was achieved]
 - Find the total number of successful and failure mission outcomes
 - Ranked the landing outcomes between 2010-06-04 to 2017-03-20
- GitHub URL of EDA with SQL notebook: https://github.com/sureshmanjoshi/testrepo/blob/main/jupyter-labs-eda-sql-coursera_sqllite.ipynb

Build an Interactive Map with Folium

- We added map objects such as markers, circles, lines to a folium map.
- We added circle for locate all launch sites in the map.
- We used marker color green for success (1) and red for failure (0) to a folium map.
- We used lines to calculate distance between launch site and selected points such as coastline, railway, highway, city.

GitHub URL of interactive map with Folium map:
 https://github.com/sureshmanjoshi/testrepo/blob/main/lab_jupyter_launch_site_location.ipynb

Build a Dashboard with Plotly Dash

- We created interactive Spacex Launch Record Dashboard.
- We plotted interactive pie chart to display success rate of launch sites. If selected all sites, it displays all sites success rate. If we select specific site, then it will plot the success and failure rate of the launch.
- We have also plotted scatter plot to display the correlation between payload and success of the launch site.
- GitHub URL of Plotly Dash lab: https://github.com/sureshmanjoshi/testrepo/blob/main/spacex_dash_app.py

Predictive Analysis (Classification)

- We used Pandas to load data into data frame. Then we manipulate some data using Numpy for required format. Then we standardized the data using StandardScalar and transform methods.
- We split the data into training and testing data. Then the models are trained and hyperparameters are selected using the GridSearchCV function.
- We used accuracy as the metric in our model. Then we used confusion matrix to evaluate the performance of our model.
- We found the classification model as best performance model.
- GitHub URL predictive analysis lab: https://github.com/sureshmanjoshi/testrepo/blob/main/SpaceX_Machine_Learning_Prediction_Part_5.ipynb

Predictive Analysis (Classification) - Flow chart

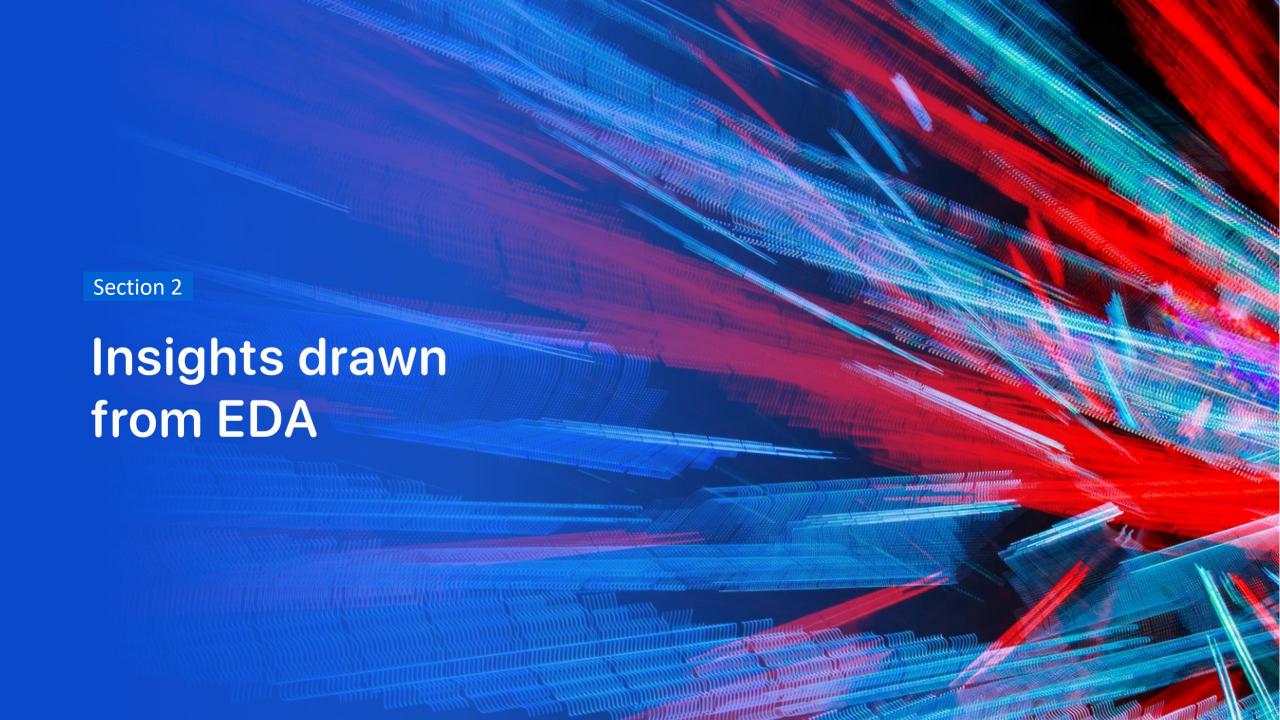
```
from js import fetch
import io
URL1 = "https://cf-courses-data.s3.us.cloud-object-stora-
respl = await fetch(URL1)
text1 = io.BytesIO((await respl.arrayBuffer()).to py())
data = pd.read csv(text1)
URL2 = 'https://cf-courses-data.s3.us.cloud-object-storag
resp2 = await fetch(URL2)
text2 = io.BytesIO((await resp2.arrayBuffer()).to py())
X = pd.read csv(text2)
Y = data['Class'].to numpy()
# students get this
transform = preprocessing.StandardScaler()
scaler = preprocessing.StandardScaler().fit(X)
X = scaler.transform(X)
X train, X test, Y train, Y test = train test split(X, Y, test size=0.2, random state=2)
parameters ={"C":[0.01,0.1,1],'penalty':['12'], 'solver':['lbfgs']}#
lr=LogisticRegression()
gscv = GridSearchCV(lr, parameters, scoring='accuracy', cv=10)
logreg cv = gscv.fit(X train, Y train)
```

```
yhat=logreg cv.predict(X test)
  plot confusion matrix(Y test, yhat)
                           Confusion Matrix
   did not land
                                                                            - 10
                                                     3
True labels
   landed
                                                    12
                did not land
                                                   land
                             Predicted labels
```

logreg cv.score(X test, Y test)

Results

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

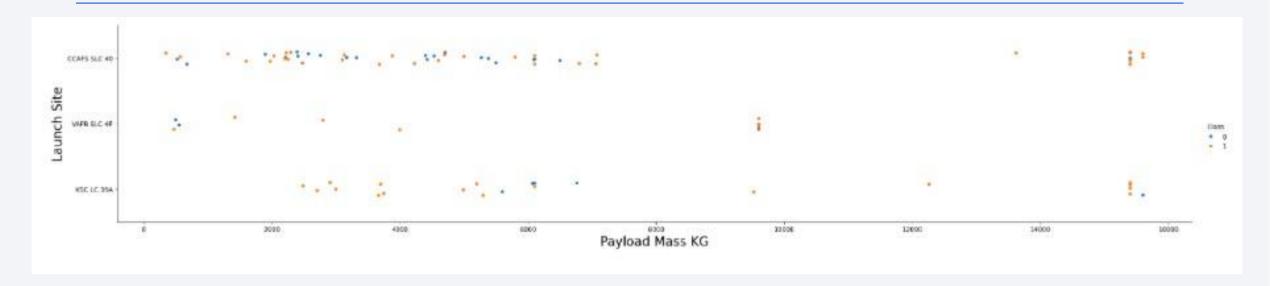


Flight Number vs. Launch Site



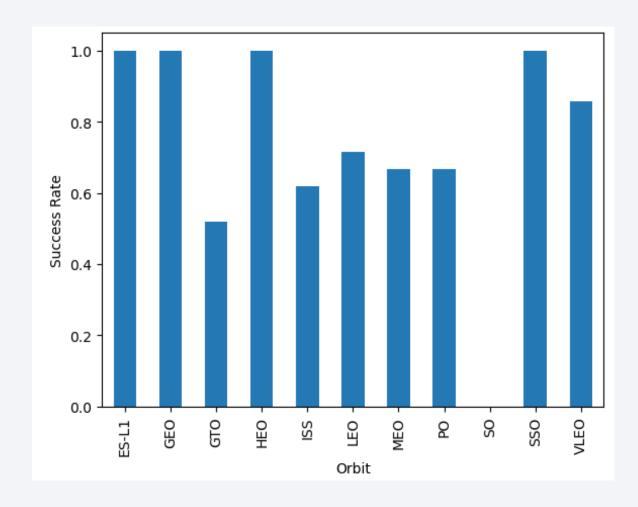
• From the above plot, we see that the greater the flight number greater the success rate of launch site.

Payload vs. Launch Site



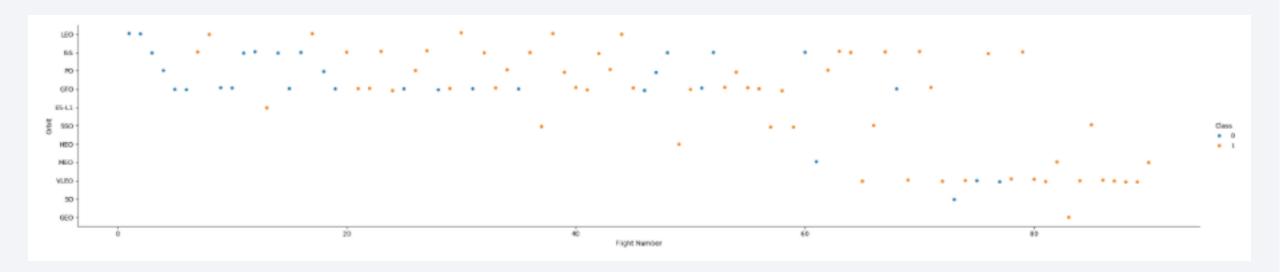
• From the above plot, we will find that for VAFB-SLC launch site there are no rockets launched for heavypayload mass > 10000

Success Rate vs. Orbit Type



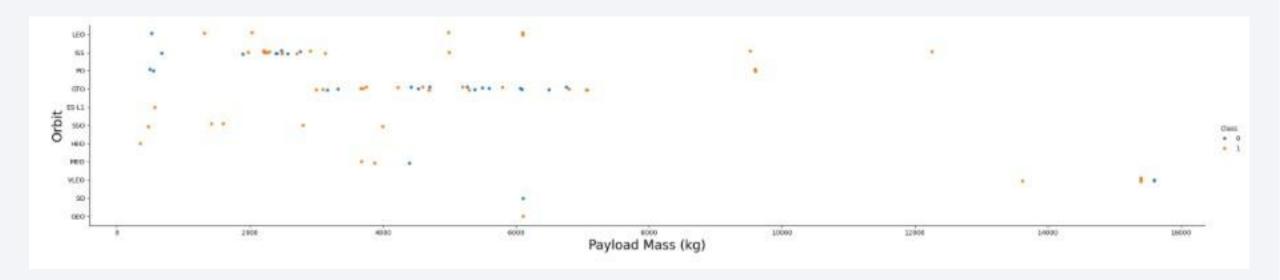
 From the Bar Plot, it is clear that ES-L1, GEO, HEO, SSO and VLEO had the most success rate.

Flight Number vs. Orbit Type



 We can see that 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.

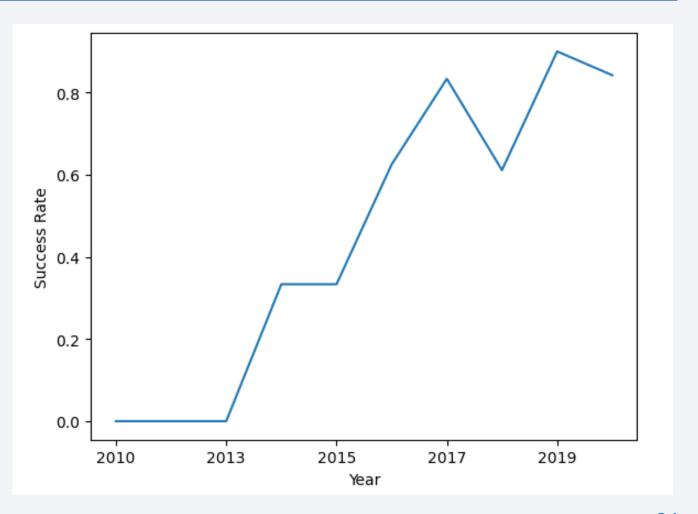
Payload vs. Orbit Type



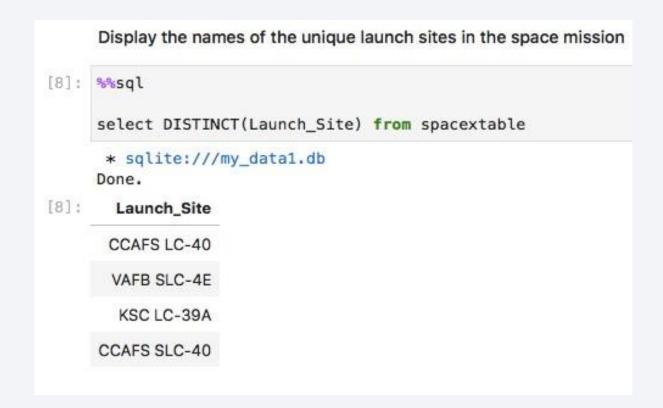
- With heavy payloads the successful landing rate are more for Polar, LEO and ISS.
- However for GTO we cannot distinguish this well as both positive and negative landing rate are there.

Launch Success Yearly Trend

• We can observe that the success rate since 2013 kept increasing till 2020.



All Launch Site Names



 We used DISTINCT keyword in SQL query to show unique launch sites of SpaceX data set.

Launch Site Names Begin with 'CCA'

We used like and limit keywords to display following result

	%sql s	elect *	from spacextabl	e where Laur	ch_Site lik	e "CCA%" limit 5			□ 个	↑ 4 7 ■	
	* 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 attemp	
	2012- 10- 08	0:35:00	F9 v1.0 B0006	CCAFS LC- 40	SpaceX CRS-1	500	LEO (ISS)	NASA (CRS)	Success	No attemp	
	2013- 03- 01	15:10:00	F9 v1.0 B0007	CCAFS LC- 40	SpaceX CRS-2	677	LEO (ISS)	NASA (CRS)	Success	No attemp	

Total Payload Mass

We used sum() method in SQL Query to calculate total payload mass.

```
Display the total payload mass carried by boosters launched by NASA (CRS)

[14]: %sql select sum(PAYLOAD_MASS__KG_) as "Total Payload Mass" from spacextable where Customer = "NASA (CRS)"

* sqlite://my_datal.db
Done.

[14]: Total Payload Mass

45596
```

Average Payload Mass by F9 v1.1

We used average() function to calculate average payload mass.

```
Display average payload mass carried by booster version F9 v1.1

[16]: %sql select avg(PAYLOAD_MASS__KG_) as "Average Payload Mass" from spacextable where Booster_Version = "F9 v1.1"

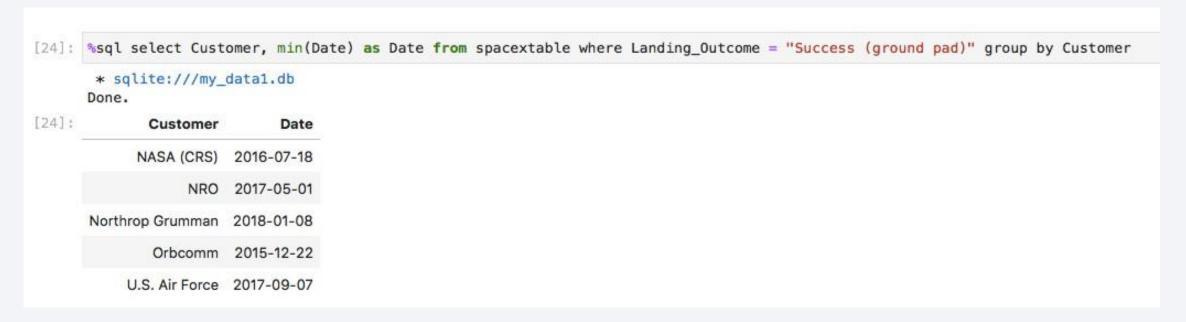
* sqlite:///my_data1.db
Done.

[16]: Average Payload Mass

2928.4
```

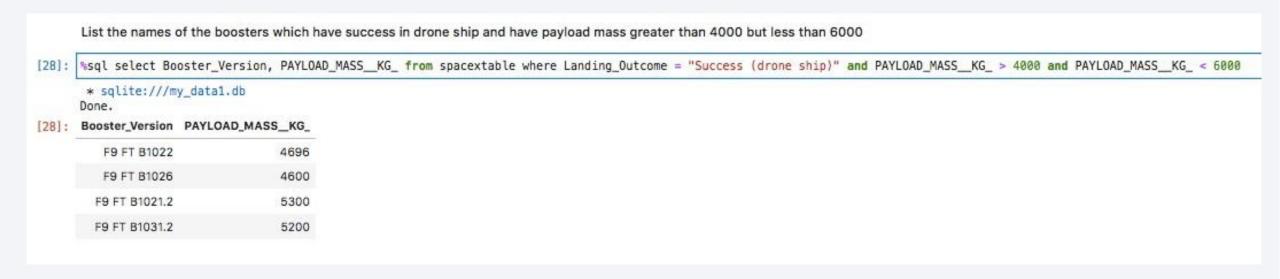
First Successful Ground Landing Date

We used min() function and applied group by keyword on customer.



Successful Drone Ship Landing with Payload between 4000 and 6000

 We used and condition to determine payload mass between 4000 to 6000 and used where clause to filter successful drone ship landing.



Total Number of Successful and Failure Mission Outcomes

We applied count() function on Mission_Outcome and group by keyword.

List the total	number of success	sful and failure m	ission outcomes						
1]: %sql select	%sql select Mission_Outcome, count(Mission_Outcome) as "Total Number" from spacextable group by Mission_Outcome								
* sqlite:/ Done.	///my_data1.db								
1]:	Mission_Outcome	Total Number							
	Failure (in flight)	1							
	Success	98							
	Success	1							
Success (pay	load status unclear)	1							

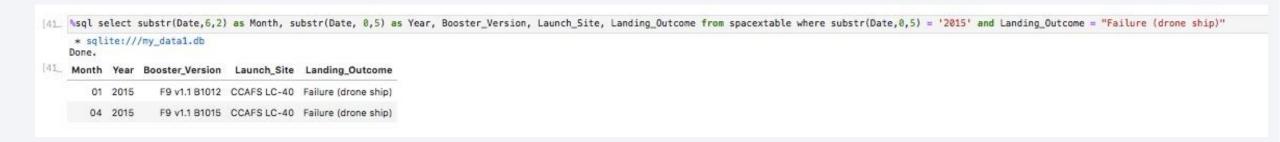
Boosters Carried Maximum Payload

 We use subquery and max() function to determine Boosters carried maximum payload.

```
List the names of the booster_versions which have carried the maximum payload mass. Use a subquery
[10]: %%sql
      select Booster_Version, PAYLOAD_MASS__KG_ as "Payload_Mass" from spacextable where PAYLOAD_MASS__KG_ = (select max(PAYLOAD_MASS__KG_) from spacextable)
       * sqlite:///my_data1.db
      Done.
      Booster_Version Payload_Mass
        F9 B5 B1048.4
                              15600
         F9 B5 B1049.4
                              15600
         F9 B5 B1051.3
                              15600
        F9 B5 B1056.4
                              15600
        F9 B5 B1048.5
                              15600
         F9 B5 B1051.4
                              15600
         F9 B5 B1049.5
                              15600
         F9 B5 B1060.2
                              15600
        F9 B5 B1058.3
                              15600
         F9 B5 B1051.6
                              15600
        F9 B5 B1060.3
                              15600
         F9 B5 B1049.7
                              15600
```

2015 Launch Records

• Sqlite does not support month names, so we used substr() function on date field to split month and date and where condition for filter.



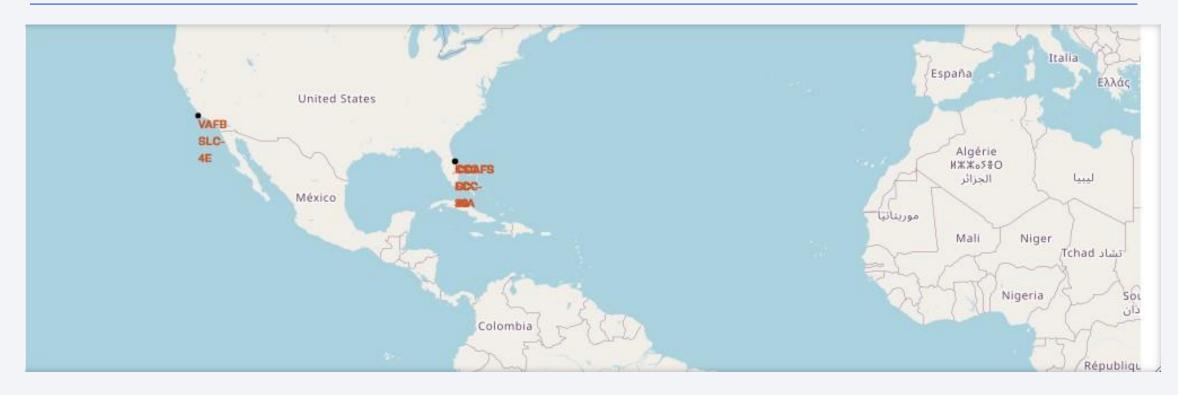
Rank Landing Outcomes Between 2010-06-04 and 2017-03-20

• We used count() function and group by keyword to determine rank and used where clause with between and and conditions to filter data.

ct Date,	Landing_Outco	ne, count(Landing_Ou	Outcome) as "Landing Outcome" from spacextable where Date between "2010-06-04" and "2017-03-20" group by Landing_Outcome order by count(Landing_Outcome) desc
lite	:///my_data1.db		
Da	te Landing_Outcome	Landing Outcome	
- 5	2 No attempt	10	
	8 Success (drone ship)	5	
10	Failure (drone ship)	5	
22 Su	ccess (ground pad)	3	
4-	8 Controlled (ocean)	3	
9-2	9 Uncontrolled (ocean)	2	
-(4 Failure (parachute)	2	
5-06-2	8 Precluded (drone ship)	1	

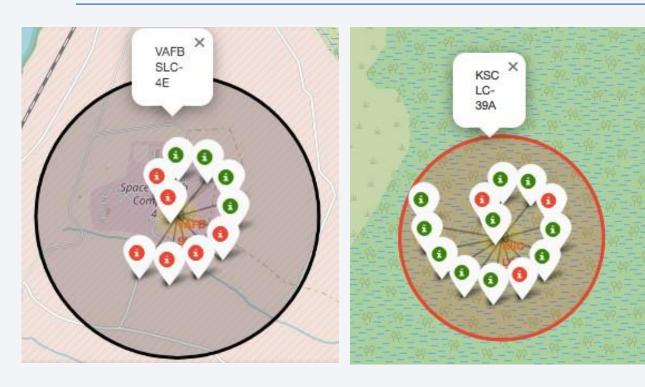


Launch Sites in global map



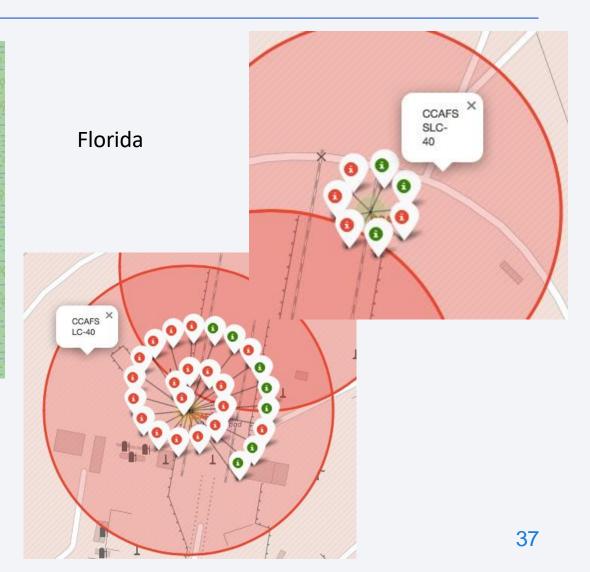
• We can see that all launch sites are in proximity to the equator line and all launch sites are very close proximity to the American coast, Florida and California.

Launch sites with color markers

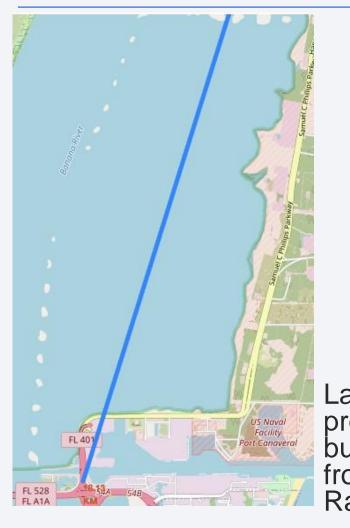


California

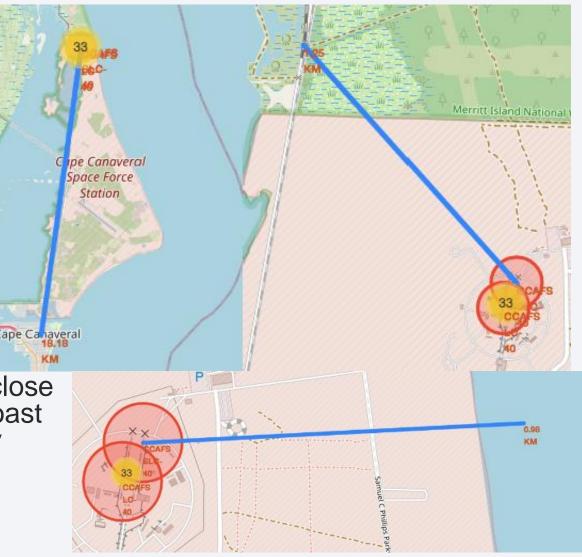
Green Marker indicates Successful launch and Red Marker indicates Failure launch.

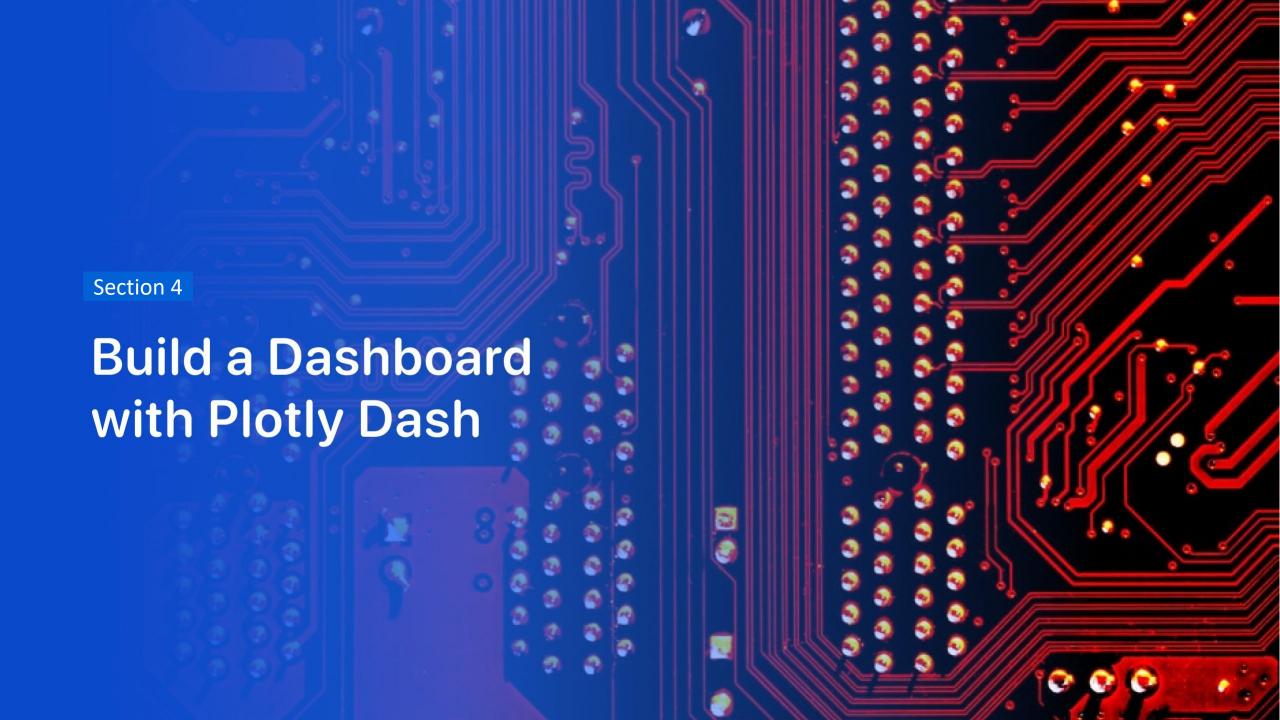


Launch site distance from Coast, Railway, Highway and City

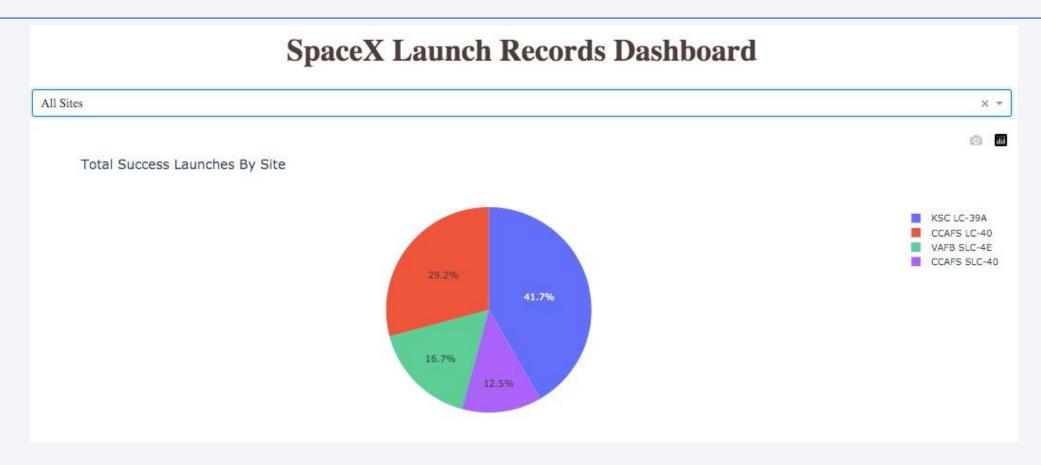


Launch sites are close proximity to the coast but distance away from Highway, Railway and City



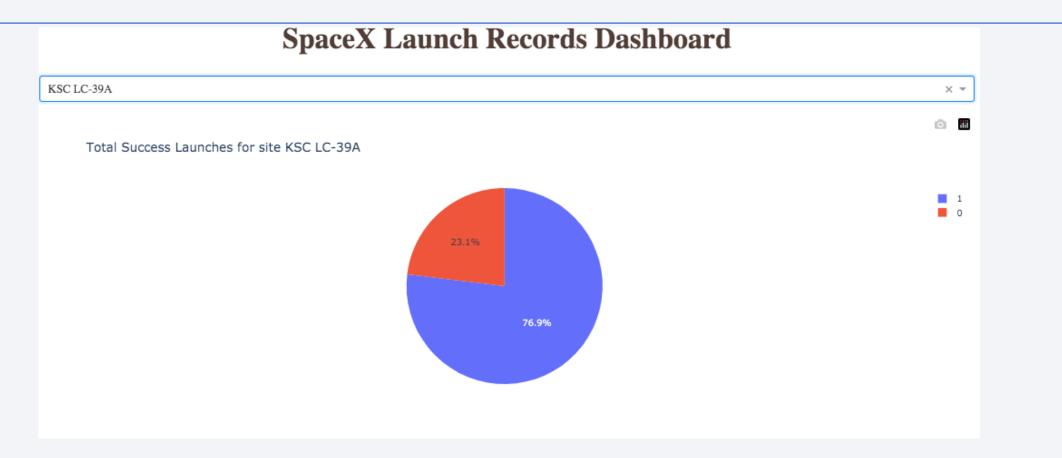


Success percentage of each launch site with Pie chart



• KSC LC - 39A has most success rate of 41.7% among all other launch site

Launch site with highest launch success rate with Pie chart



• KSC LC - 39A has 76.9% success and 23.1% failure rate, which is highest success rate among other three launch sites.

41

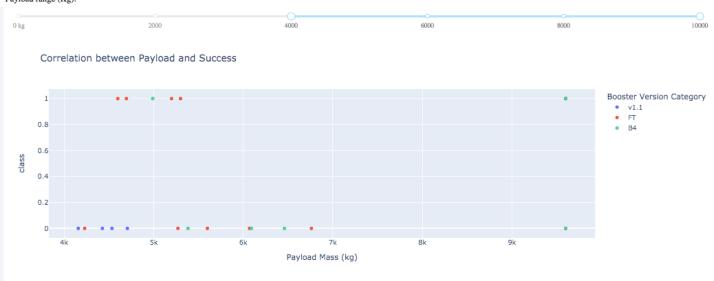
Payload vs. Launch Outcome scatter plot



0kg to 4000 kg payload

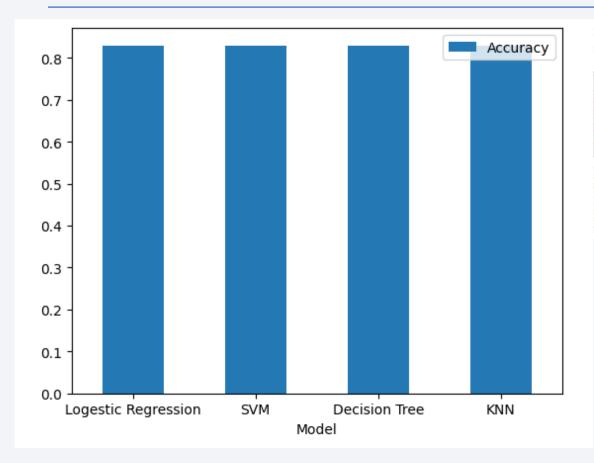
4000 kg to 10000 kg payload

- It is clearly seen that low weighted payload(0-4000kg) has most success rate than heavy weighted(4000-10000kg) payload
- FT Booster Version has the highest launch success rate.





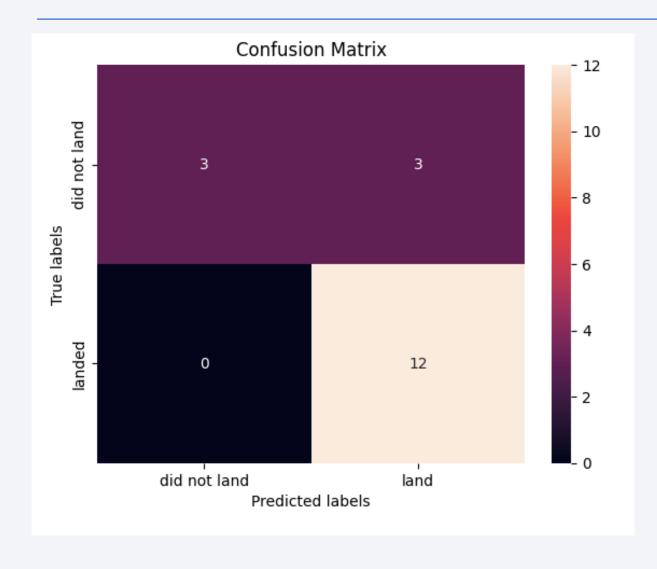
Classification Accuracy



Find the method performs best:

They all perform practically the same.

Confusion Matrix



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

Conclusions

- The greater the flight number greater the success rate of launch site.
- ES-L1, GEO, HEO, SSO and VLEO had the most success rate.
- With heavy payloads the successful landing rate are more for Polar, LEO and ISS.
- The success rate since 2013 kept increasing till 2020.
- All launch sites are in proximity to the equator line and all launch sites are very close proximity to the American coast, Florida and California
- Launch sites are close proximity to the coast but distance away from Highway, Railway and City
- Low weighted payload (0-4000kg) has most success rate than heavy weighted (4000-10000kg) payload
- KSC LC 39A has 76.9% success and 23.1% failure rate, which is highest success rate among other three launch sites.

Appendix

- The GitHub URL of the SpaceX API calls notebook: https://github.com/sureshmanjoshi/testrepo/blob/main/jupyter-labs-spacex-data-collection-api.ipynb
- GitHub URL of the web scraping notebook: https://github.com/sureshmanjoshi/testrepo/blob/main/jupyter-labs-webscraping.ipynb
- GitHub URL of data wrangling related notebooks: <u>https://github.com/sureshmanjoshi/testrepo/blob/main/labs-jupyter-spacex-Data wrangling.ipynb</u>
- GitHub URL of EDA with data visualization notebook: https://github.com/sureshmanjoshi/testrepo/blob/main/jupyter-labs-eda-dataviz.ipynb
- GitHub URL of EDA with SQL notebook: https://github.com/sureshmanjoshi/testrepo/blob/main/jupyter-labs-eda-sql-coursera_sqllite.ipynb
- GitHub URL of interactive map with Folium map: https://github.com/sureshmanjoshi/testrepo/blob/main/lab_jupyter_launch_site_location.ipynb
- GitHub URL of Plotly Dash lab: https://github.com/sureshmanjoshi/testrepo/blob/main/spacex_dash_app.py
- GitHub URL predictive analysis lab: https://github.com/sureshmanjoshi/testrepo/blob/main/SpaceX_Machine_Learning_Prediction_Part_5.ipynb

