



IBM Developer
SKILLS NETWORK

Winning Space Race with Data Science

Franco Messetzián
October 2021



OUTLINE

- EXECUTIVE SUMMARY
- INTRODUCTION
- METHODOLOGY
- RESULTS
- CONCLUSION
- APPENDIX

EXECUTIVE SUMMARY

- SUMMARY OF METHODOLOGIES
 - DATA COLLECTION THROUGH API
 - DATA COLLECTION WITH WEB SCRAPING
 - DATA WRANGLING
 - EXPLORATORY DATA ANALYSIS WITH SQL
 - EXPLORATORY DATA ANALYSIS WITH DATA VISUALIZATION
 - INTERACTIVE VISUAL ANALYTICS WITH FOLIUM
 - MACHINE LEARNING PREDICTION
- SUMMARY OF ALL RESULTS
 - EXPLORATORY DATA ANALYSIS RESULT
 - INTERACTIVE ANALYTICS SHOWN IN SCREENSHOTS
 - PREDICTIVE ANALYTICS RE_{SULT}

INTRODUCTION

- SPACEX IS A COMPANY THAT DESIGNS, MANUFACTURES AND LAUNCHES ADVANCED ROCKETS AND SPACECRAFT. ONE OF THEIR MAIN COMPETITIVE ADVANTAGES IN THE ROCKET LAUNCHING MARKET IS THE COST AT WHICH THEY OFFER THEIR LAUNCH, WHICH IS NEARLY A THIRD OF THE PRICE OF THEIR COMPETITORS. THEY ACHIEVE THIS BY RE USING THE FIRST STAGE.
- IN THIS PROJECT WE ARE GOING TO ANALYZE THE AVAILABLE INFORMATION IN ORDER TO DETERMINE (AND THEREFORE BE ABLE TO CONFIDENTLY PREDICT) THE BEST INDICATORS FOR THEIR FALCON9 SUCCESSFUL LANDINGS. IN TERM, THIS WILL ACCURATELY DETERMINE THE COST OF EACH LAUNCH, AND VALUABLE INFORMATION FOR A SPACEX COMPETITOR.

Section 1

Methodology

METHODOLOGY

- EXECUTIVE SUMMARY
- DATA COLLECTION METHODOLOGY:
 - DATA WAS COLLECTED BY WEB SCRAPING SPACEX API'S 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
 - HOW TO BUILD, TUNE, EVALUATE CLASSIFICATION MODELS

DATA COLLECTION



DATA COLLECTION – SPACEX API

- WE USED THE GET REQUEST COMMAND ON SPACEX API TO COLLECT AND CLEAN THE DATA. THEN DID SOME BASIC DATA WRANGLING AND FORMATTING.
- CLICK ON THE [DATA COLLECTION NOTEBOOK](#) FOR DETAILS

Task 1: Request and parse the SpaceX launch data using the GET request

To make the requested JSON results more consistent, we will use the following static response object for this project:

```
11]: static_json_url='https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBM-D50321EN-SkillsNetwork/datasets/API_call_spacex_ap:
```

We should see that the request was successful with the 200 status response code

```
12]: response.status_code
```

```
12]: 200
```

Now we decode the response content as a Json using `.json()` and turn it into a Pandas dataframe using `.json_normalize()`

```
13]: # Use json_normalize method to convert the json result into a dataframe
data = pd.json_normalize(response.json())
```

DATA COLLECTION - SCRAPING

- WE DID WEB SCRAPING FOR FALCON9 LAUNCH RECORDS USING BEAUTIFULSOUP
- WE PARSED THE TABLE AND CONVERTED IT INTO A PANDAS DATAFRAME
- CLICK ON THE WEB SCRAPING NOTEBOOK FOR DETAILS

TASK 1: Request the Falcon9 Launch Wiki page from its URL

First, let's perform an HTTP GET method to request the Falcon9 Launch HTML page, as an HTTP response.

```
[5]: # use requests.get() method with the provided static_url  
# assign the response to a object  
page = requests.get(static_url).text
```

Create a `BeautifulSoup` object from the HTML `response`

```
[6]: # Use BeautifulSoup() to create a BeautifulSoup object from a response text content  
soup = BeautifulSoup(page, "html.parser")
```

Print the page title to verify if the `BeautifulSoup` object was created properly

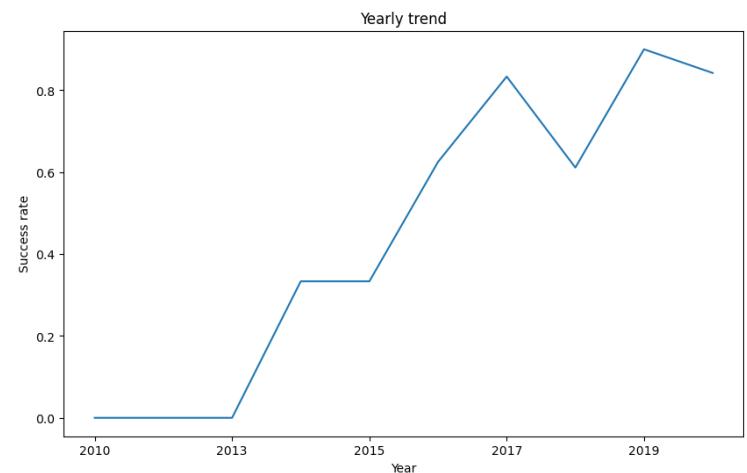
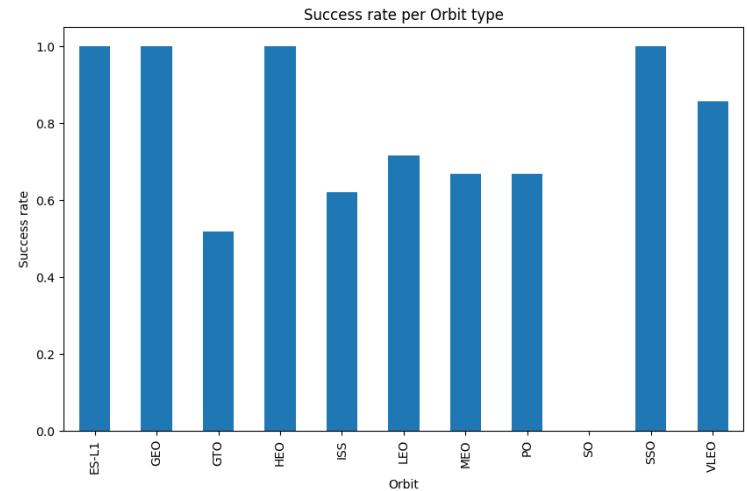
```
[9]: # Use soup.title attribute  
print(soup.title)  
  
<title>List of Falcon 9 and Falcon Heavy launches - Wikipedia</title>
```

DATA WRANGLING

- WE PERFORMED EXPLORATORY DATA ANALYSIS AND DETERMINED TRAINING LEVELS
- WE CALCULATED THE NUMBER OF LAUNCHES AT EACH SITE AND THE NUMBER OF OCCURANCE OF EACH ORBIT
- WE THEN CREATED A LANDING OUTCOME LABEL FROM THE OUTCOME COLUMN OF OUR DATAFRAME
- REFER TO THE DATA WRANGLING NOTEBOOK FOR DETAILS

BUILD A DASHBOARD WITH PLOTLY DASH

- WE EXPLORED THE DATA BY VISUALIZING THE RELATIONSHIP BETWEEN DIFFERENT VARIABLES IN ORDER TO DETECT DEPENDENCIES FOR FURTHER ANALYSIS.
- THE GRAPHS ON THE RIGHT SHOW HOW CERTAIN ORBITS HAVE BETTER SUCCESS RATE (TOP), AND THE GENERAL POSITIVE RELATIONSHIP BETWEEN SUCCESS RATE AND YEARS SINCE FIRST LAUNCH (BOTTOM)
- SEE [EDA DATA VISUALIZATION NOTEBOOK](#) FOR DETAILS



EDA WITH SQL

- WE LOADED THE SPACEX DATASET INTO A POSTGRESQL DATABASE
- USING SQL COMMANDS, WE EXPLORED THE DATASET FOR INSIGHTS SUCH US:
 - NAMES OF UNIQUE LAUNCH SITES
 - TOTAL PAYLOAD MASS CARRIED BY BOOSTERS LAUNCHED BY NASA (CRS)
 - AVERAGE PAYLOAD MASS CARRIED BY BOOSTER VERSION F9 v1.1
 - FAILED LANDING OUTCOMES IN DRONE SHIP WITH THEIR BOOSTER VERSION AND LAUNCH SITE NAMES
- SEE [EDA WITH SQL](#) FOR FURTHER DETAILS

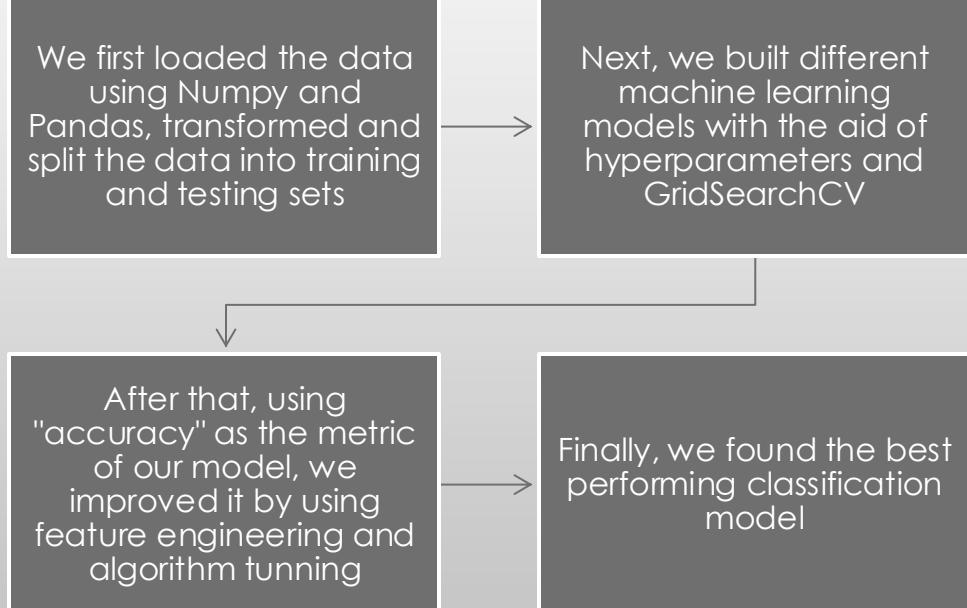
BUILD AN INTERACTIVE MAP WITH FOLIUM

- WE CREATED A FOLIUM MAP, THEN MARKED ALL LAUNCH SITES AND ADDED MAP OBJECTS SUCH AS MARKERS, CIRCLES AND LINES.
- THESE MARKERS HELPED US VISUALIZE THE SUCCESS/FAILURE LAUNCHES FOR EACH SITE ON THE MAP
- REFER TO FOLIUM NOTEBOOK FOR FULL DETAILS

BUILD A DASHBOARD WITH PLOTLY DASH

- WE BUILT AN INTERACTIVE DASHBOARD WITH PLOTLY DASH
- WE PLOTTED PIE CHARTS SHOWING THE TOTAL LAUNCHES GIVEN A SITE
- WE ALSO PLOTTED A SCATTER GRAPH FOR DIFFERENT BOOSTER VERSIONS, WHERE WE CAN SEE PAYLOAD AND OUTCOME FOR EACH
- SEE [DASHBOARD NOTEBOOK](#) FOR DETAILS

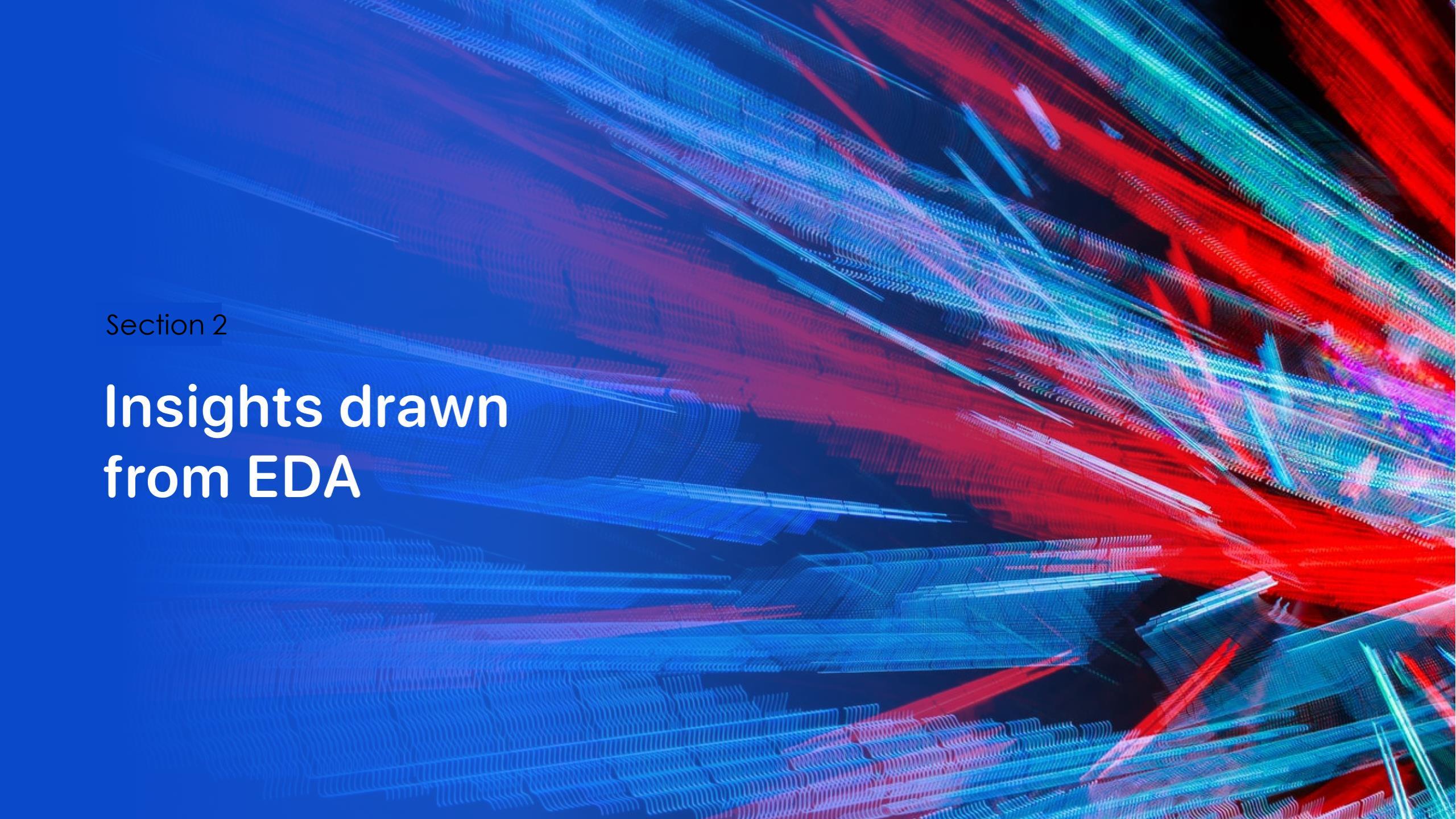
PREDICTIVE ANALYSIS (CLASSIFICATION)



- FIND THE FULL DETAILS AT [PREDICTIVE ANALYSIS NOTEBOOK](#)

RESULTS

- EXPLORATORY DATA ANALYSIS RESULTS
- INTERACTIVE ANALYTICS DEMO IN SCREENSHOTS
- PREDICTIVE ANALYSIS RESULTS

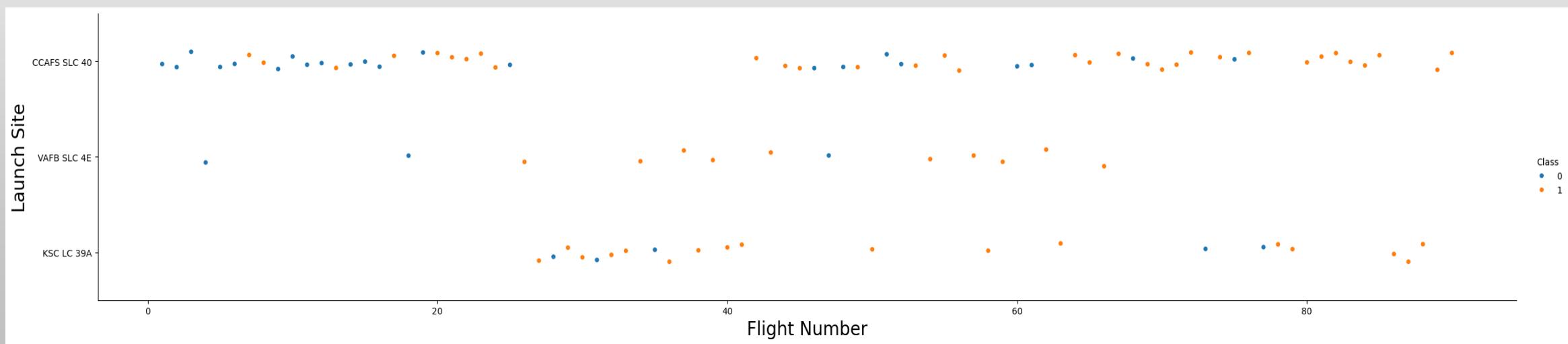
The background of the slide features a complex, abstract digital visualization. It consists of a grid of points that have been connected by thin lines, creating a three-dimensional effect. The colors used are primarily shades of blue, red, and green, with some purple and yellow highlights. The overall appearance is reminiscent of a microscopic view of a crystal lattice or a complex data visualization.

Section 2

Insights drawn from EDA

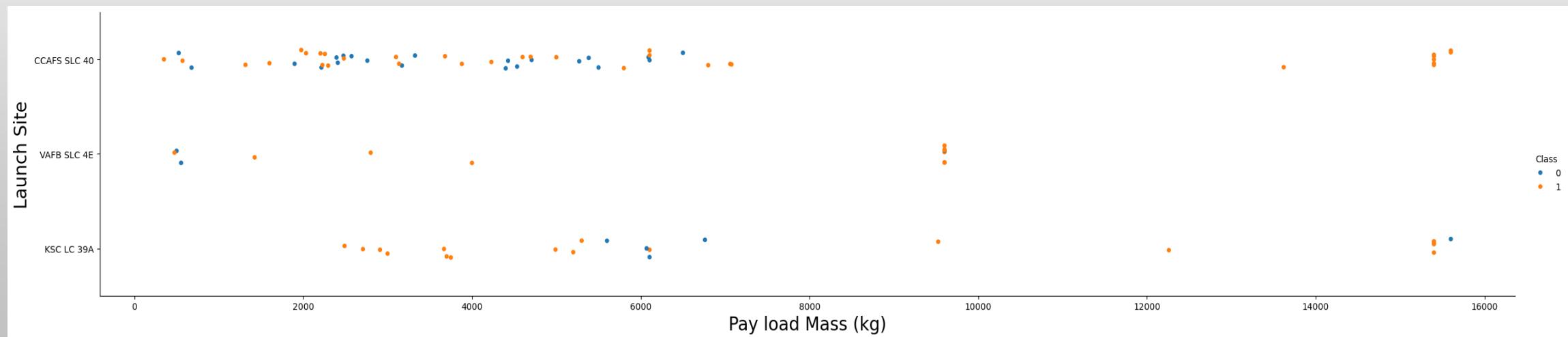
FLIGHT NUMBER VS. LAUNCH SITE

- WE CAN SEE ON THE PLOT THAT AS FLIGHT NUMBERS INCREASED, SO DID THE SUCCESS RATE FOR EACH SITE



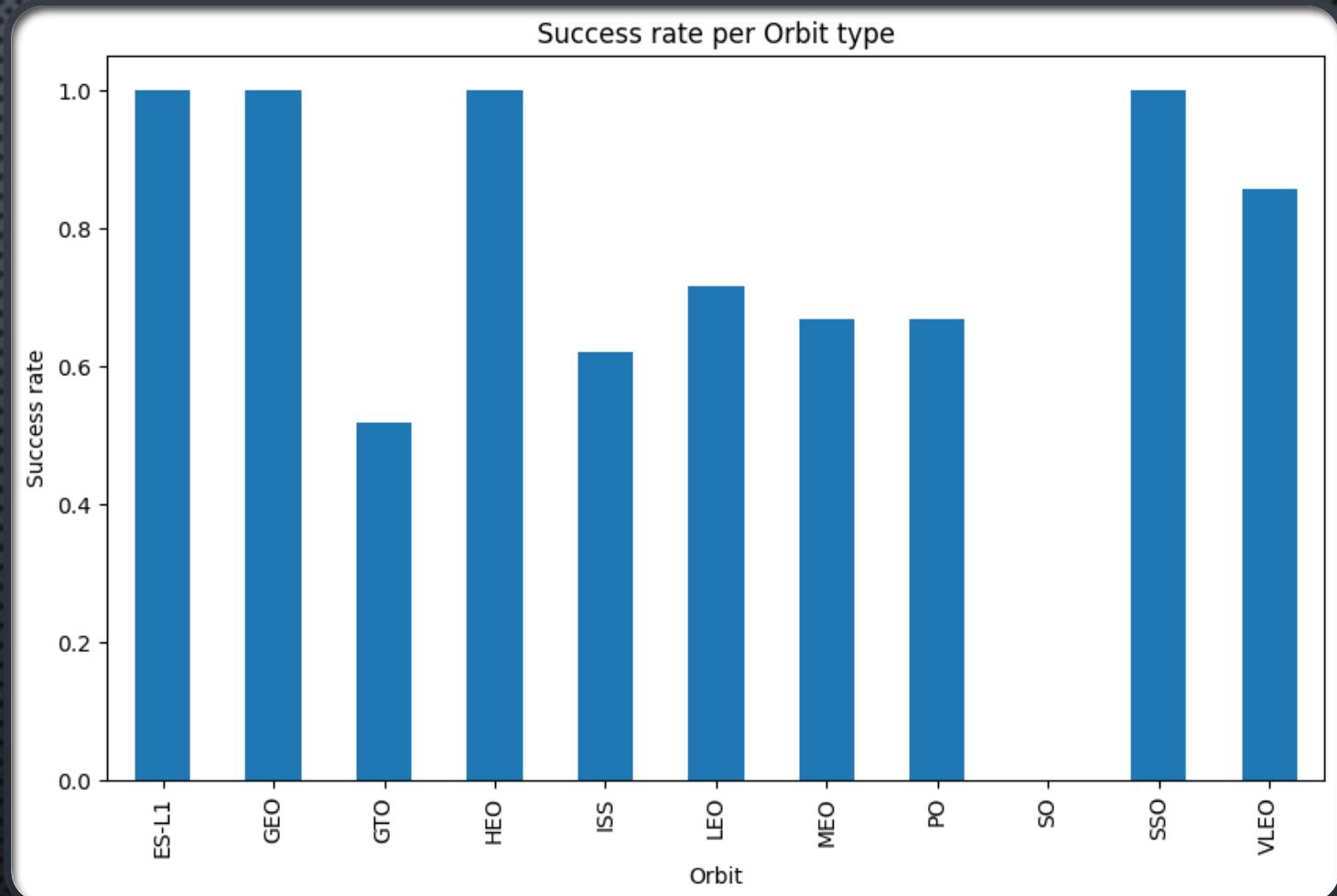
PAYLOAD VS. LAUNCH SITE

- WE CAN OBSERVE THAT FOR THE VAFB-SLC LAUNCH SITE THERE ARE NO ROCKETS LAUNCHED FOR HEAVY PAYLOAD MASS(GREATER THAN 10000)



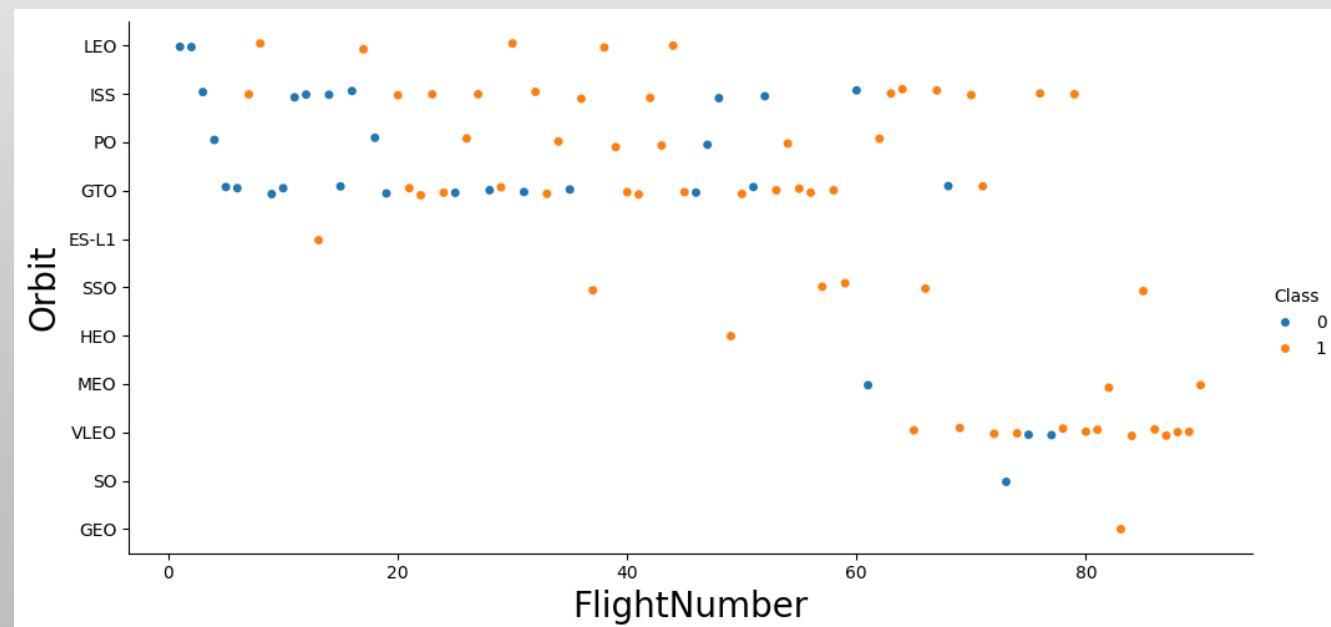
SUCCESS RATE VS. ORBIT TYPE

- WE CAN CLEARLY ESTABLISH THAT ES-L1, GEO, HEO AND SSO LEAD THE ORBIT TYPE SUCCESS RATE



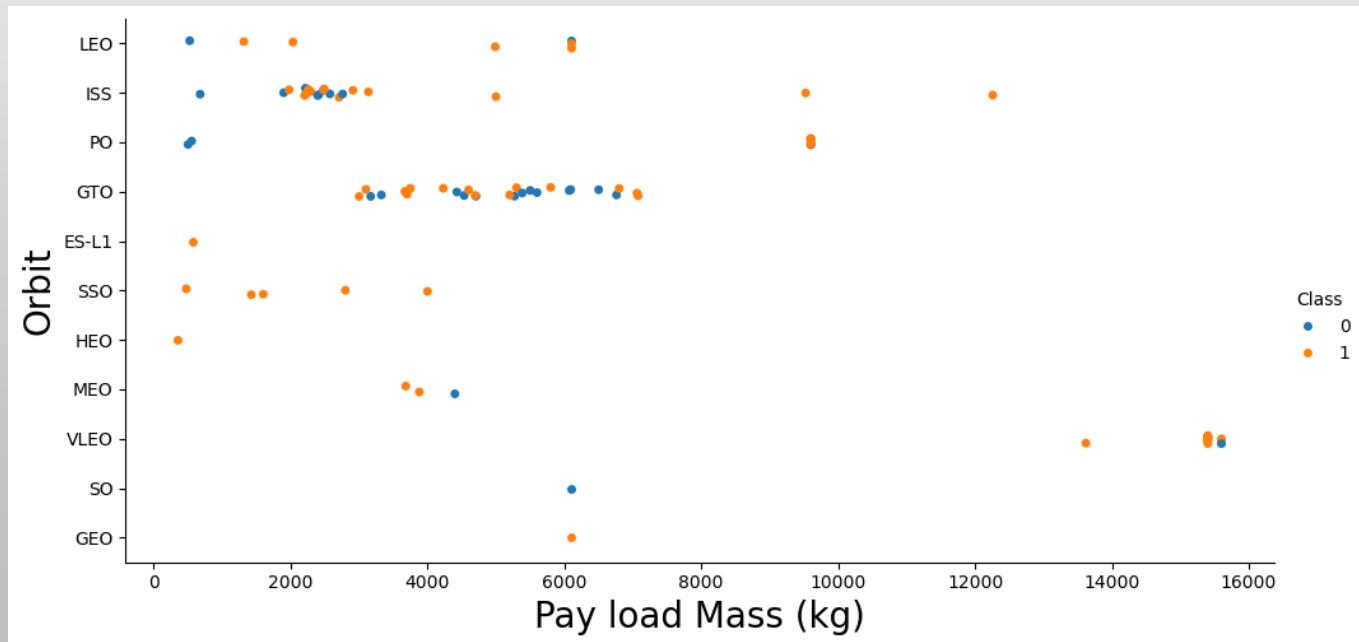
FLIGHT VS. ORBIT TYPE

- IT'S INTERESTING TO NOTE THAT FOR SOME ORBITS, SUCH AS LEO, THERE IS A DIRECT RELATIONSHIP BETWEEN THE NUMBER OF FLIGHTS AND SUCCESS RATE. ON THE CONTRARY, ORBITS LIKE GTO OR ISS, THAT DOESN'T SEEM THE CASE



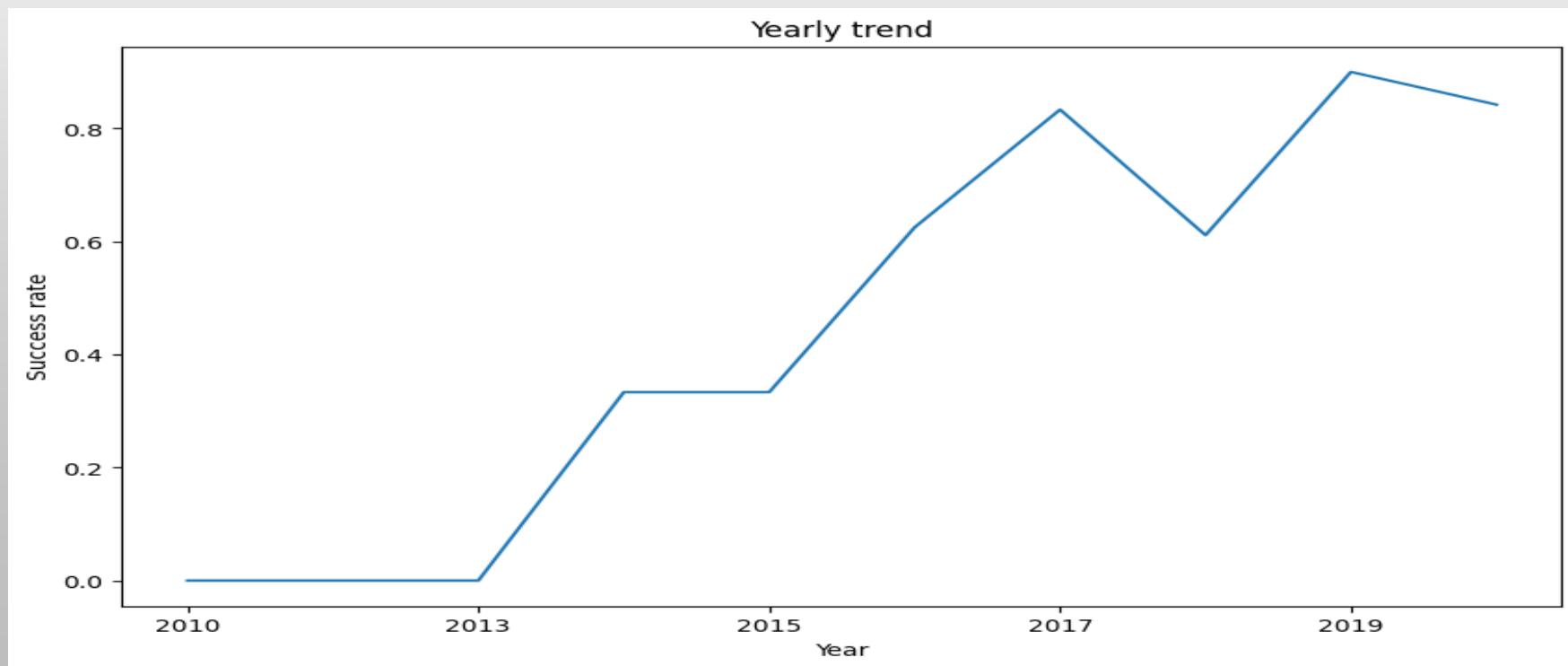
PAYOUT VS. ORBIT TYPE

- WE CAN SEE THAT FOR HEAVY PAYLOADS, SUCCESS RATE CORRELATES MORE WITH PO, LEO AND ISS ORBITS



LAUNCH SUCCESS YEARLY TREND

- FROM THE GRAPH WE CAN SEE HOW SUCCESS INCREASED STEADILY FROM 2013 ONWARDS



ALL LAUNCH SITE NAMES

- USING THE COMMAND DISTINCT IN SQL WE CAN GET THE SITES LIST

Display the names of the unique launch sites in the space mission

```
[7]: %sql select distinct "Launch_Site" from SPACEXTABLE  
* sqlite:///my_data1.db  
Done.
```

Launch_Site
CCAFS LC-40
VAFB SLC-4E
KSC LC-39A
CCAFS SLC-40

LAUNCH SITE NAMES BEGIN WITH 'KSC'

- USING THE CONDITIONAL 'WHERE' AND 'LIKE' IN SQL WE GET THE DESIRED OUTCOME

```
I Display 5 records where launch sites begin with the string 'KSC'

[9]: %sql select * from SPACEXTABLE where "Launch_Site" like 'KSC%' limit 5
      * sqlite:///my_data1.db
      Done.

[9]: 

| Date       | Time (UTC) | Booster_Version | Launch_Site | Payload       | PAYLOAD_MASS_KG_ | Orbit     | Customer   | Mission_Outcome | Landing_Outcome      |
|------------|------------|-----------------|-------------|---------------|------------------|-----------|------------|-----------------|----------------------|
| 2017-02-19 | 14:39:00   | F9 FT B1031.1   | KSC LC-39A  | SpaceX CRS-10 | 2490             | LEO (ISS) | NASA (CRS) | Success         | Success (ground pad) |
| 2017-03-16 | 6:00:00    | F9 FT B1030     | KSC LC-39A  | EchoStar 23   | 5600             | GTO       | EchoStar   | Success         | No attempt           |
| 2017-03-30 | 22:27:00   | F9 FT B1021.2   | KSC LC-39A  | SES-10        | 5300             | GTO       | SES        | Success         | Success (drone ship) |
| 2017-05-01 | 11:15:00   | F9 FT B1032.1   | KSC LC-39A  | NROL-76       | 5300             | LEO       | NRO        | Success         | Success (ground pad) |
| 2017-05-15 | 23:21:00   | F9 FT B1034     | KSC LC-39A  | Inmarsat-5 F4 | 6070             | GTO       | Inmarsat   | Success         | No attempt           |


```

TOTAL PAYLOAD MASS

- USING SQL's FUNCTION 'SUM' AND CONDITIONAL '

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

```
[15]: %sql select sum("PAYLOAD_MASS_KG_") from SPACEXTABLE where "Customer" = 'NASA (CRS)'  
* sqlite:///my_data1.db  
Done.  
[15]: sum("PAYLOAD_MASS_KG_")  
-----  
45596
```

AVERAGE PAYLOAD MASS BY F9 V1.1

- WE CALCULATED THE AVERAGE PAYLOAD MASS CARRIED BY F9 v1.1 AS 2928.4

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

```
16]: %sql select avg("PAYLOAD_MASS__KG_") from SPACEXTABLE where "Booster_Version" = 'F9 v1.1'  
* sqlite:///my_data1.db  
Done.  
16]: avg("PAYLOAD_MASS__KG_")  
2928.4
```

FIRST SUCCESSFUL (DRONE SHIP) GROUND LANDING DATE

- WE CAN SEE THAT THE FIRST SUCCESSFUL LANDING WAS ACHIEVED ON THE 8TH APRIL 2016

List the date where the successful landing outcome in drone ship was achieved.

Hint: Use min function

```
[36]: %sql select min("Date") from SPACEXTABLE where "Landing_Outcome" = 'Success (drone ship)'
```

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

```
Done.
```

```
[36]: min("Date")
```

```
2016-04-08
```

SUCCESSFUL DRONE SHIP LANDING WITH PAYLOAD BETWEEN 4000 AND 6000

- AS SHOWN BELOW, THE FOLLOWING BOOSTER VERSIONS HAVE SUCCESSFULLY LANDED PAYLOADS BETWEEN 4000 AND 6000: F9 FT B1032.1, F9 B4 B1040.1 AND F9 B4 B1043.1

```
List the names of the boosters which have success in ground pad and have payload mass greater than 4000 but less than 6000

39]: %sql select distinct("Booster_Version") from SPACEXTABLE where "Landing_Outcome" = 'Success (ground pad)' and "PAYLOAD_MASS_KG_" >4000 and "PAYLOAD_MASS_KG_" <6000
* sqlite:///my_data1.db
Done.

39]: Booster_Version
F9 FT B1032.1
F9 B4 B1040.1
F9 B4 B1043.1
```

TOTAL NUMBER OF SUCCESSFUL AND FAILURE MISSION OUTCOMES

- AS SHOWN BELOW, THERE HAVE BEEN 100 SUCCESSES AND 1 FAILURE

List the total number of successful and failure mission outcomes

```
: %sql select count("Mission_Outcome") as "Failure_count" from SPACEXTABLE where "Mission_Outcome" like 'Failure%'  
* sqlite:///my_data1.db  
Done.  
: Failure_count  
-----  
1  
  
: %sql select count("Mission_Outcome") as "Success_count" from SPACEXTABLE where "Mission_Outcome" like 'Success%'  
* sqlite:///my_data1.db  
Done.  
: Success_count  
-----  
100
```

BOOSTERS CARRIED MAXIMUM PAYLOAD

List the names of the booster_versions which have carried the maximum payload mass. Use a subquery

```
56]: %sql select distinct("Booster_Version") from SPACEXTABLE where "PAYLOAD_MASS_KG_" = (select max("PAYLOAD_MASS_KG_") from SPACEXTABLE)
```

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

```
Done.
```

```
56]: 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
```

- SELECTING THE DISTINCT BOOSTER VERSION AND APPLYING THE CONDITION FOR MAX PAYLOAD, WE GET THE RESULTS AS SHOWN ON THE LEFT

2017 LAUNCH RECORDS

List the records which will display the month names, successful landing_outcomes in ground pad ,booster versions, launch_site for the months in year 2017

Note: SQLlite does not support monthnames. So you need to use substr(Date,6,2) for month, substr(Date,9,2) for date, substr(Date,0,5),='2017' for year.

```
[60]: %sql select substr("Date",6,2), "Landing_Outcome", "Booster_Version", "Launch_Site" from SPACEXTABLE where "Landing_Outcome" = 'Success' (g
* sqlite:///my_data1.db
Done.

[60]: 

| substr("Date",6,2) | Landing_Outcome      | Booster_Version | Launch_Site  |
|--------------------|----------------------|-----------------|--------------|
| 02                 | Success (ground pad) | F9 FT B1031.1   | KSC LC-39A   |
| 05                 | Success (ground pad) | F9 FT B1032.1   | KSC LC-39A   |
| 06                 | Success (ground pad) | F9 FT B1035.1   | KSC LC-39A   |
| 08                 | Success (ground pad) | F9 B4 B1039.1   | KSC LC-39A   |
| 09                 | Success (ground pad) | F9 B4 B1040.1   | KSC LC-39A   |
| 12                 | Success (ground pad) | F9 FT B1035.2   | CCAFS SLC-40 |


```

- **SELECTING THE DESIRED COLUMNS, AND APPLYING THE CONDITIONS FOR SUCCESSFUL LANDING AND SUBSTR("DATE",0,5) = '2017'**
WE CAN SEE THE LAUNCH RECORDS FOR EVERY ACTIVE MONTH IN 2017

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 "Date", "Landing_Outcome", count("Landing_Outcome") as "Count of Landing Outcomes" from SPACEXTABLE where "Date" between "2010-06-04" and "2017-03-20" group by "Date", "Landing_Outcome" order by "Count of Landing Outcomes" desc;
```

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

Date	Landing_Outcome	Count of Landing Outcomes
2012-05-22	No attempt	10
2016-04-08	Success (drone ship)	5
2015-01-10	Failure (drone ship)	5
2015-12-22	Success (ground pad)	3
2014-04-18	Controlled (ocean)	3
2013-09-29	Uncontrolled (ocean)	2
2010-06-04	Failure (parachute)	2
2015-06-28	Precluded (drone ship)	1

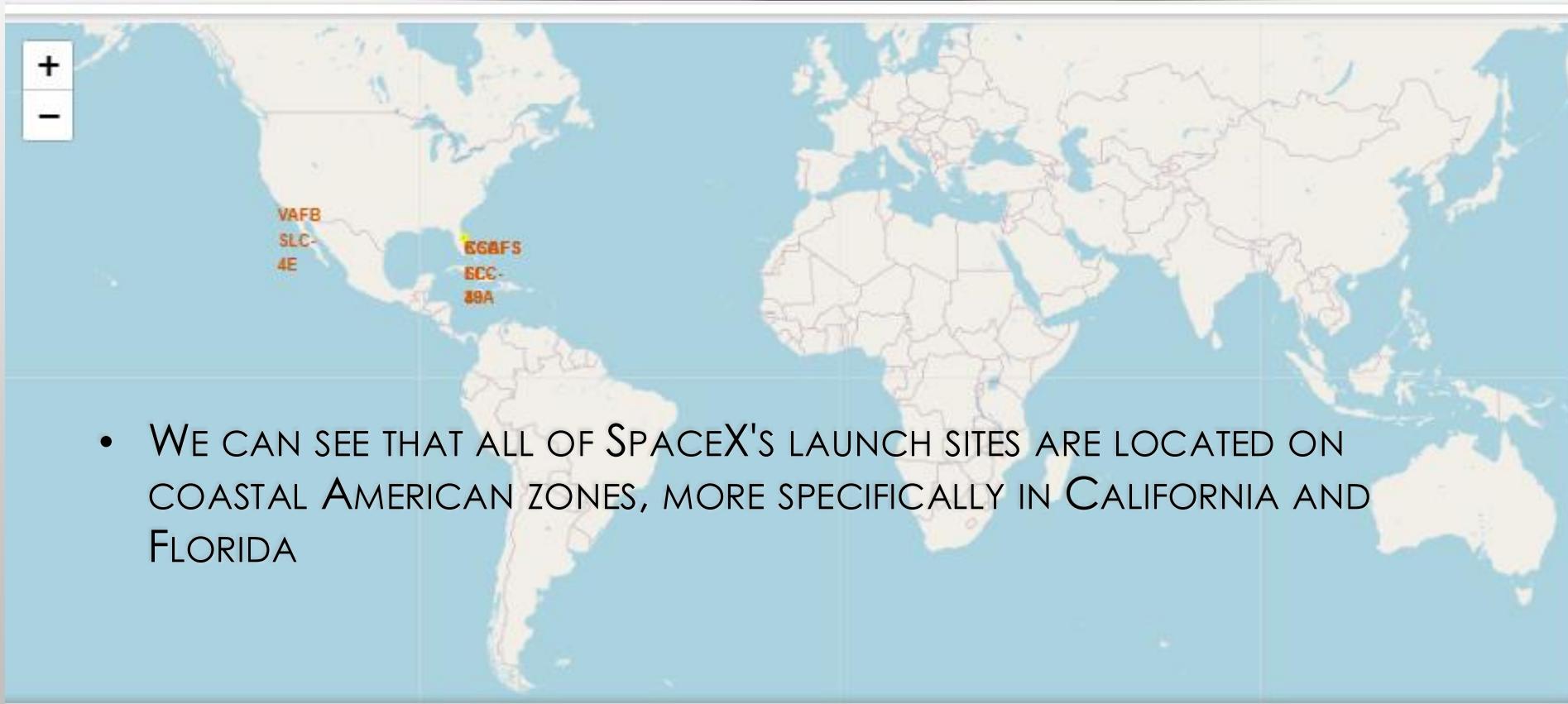
- WE SELECTED "DATE" AND COUNT OF LANDING OUTCOMES. THEN USED THE WHERE CLAUSE TO FILTER THE DESIRED RANGE.
- WE APPLIED THE GROUPBY CLAUSE ON LANDING OUTCOMES AND ORDERBY THE NEWLY CREATED COLUMN 3 IN DESCENDING ORDER

The background of the slide is a photograph taken from space at night. It shows the curvature of the Earth's horizon against a dark blue sky. City lights are visible as numerous small white and yellow dots, primarily concentrated in the lower right quadrant where the United States appears. In the upper right, there are bright green and yellow bands of light, likely the Aurora Borealis or Australis. The overall atmosphere is dark and mysterious.

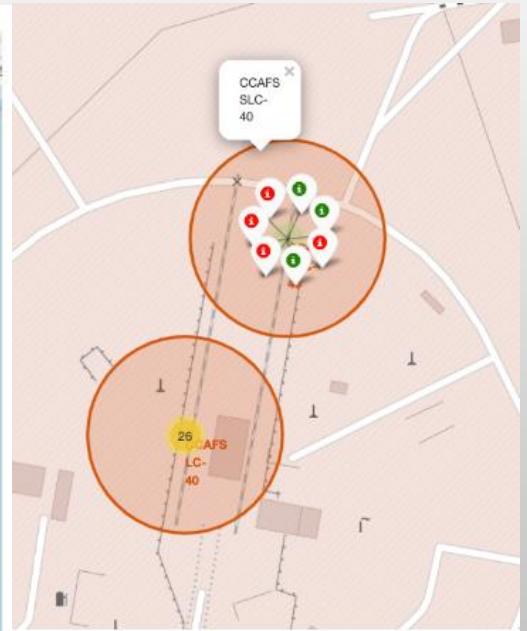
Section 3

Launch Sites Proximities Analysis

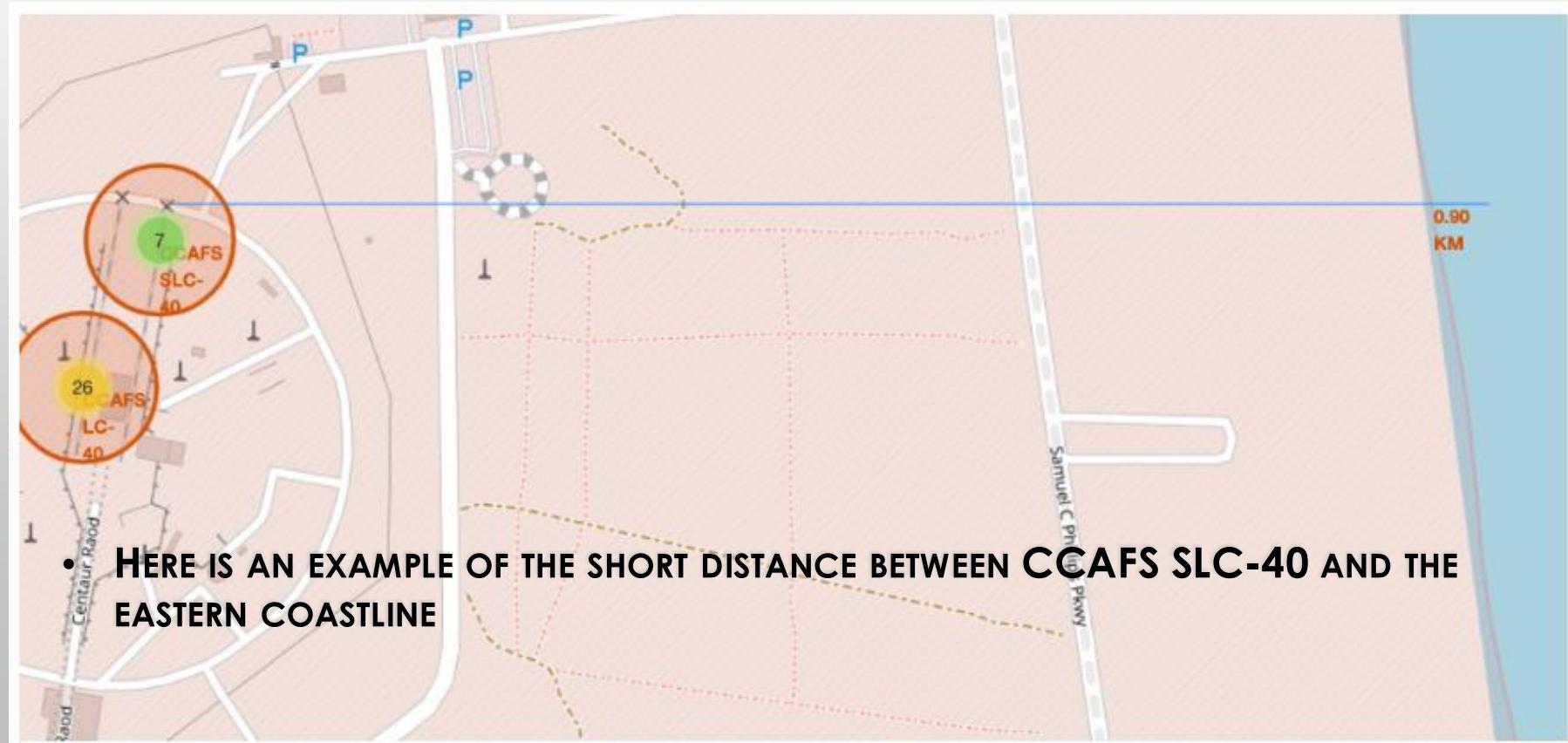
ALL LAUNCH SITES WORLD MAP



MARKERS FOR COLOR LABELED LANDING OUTCOMES

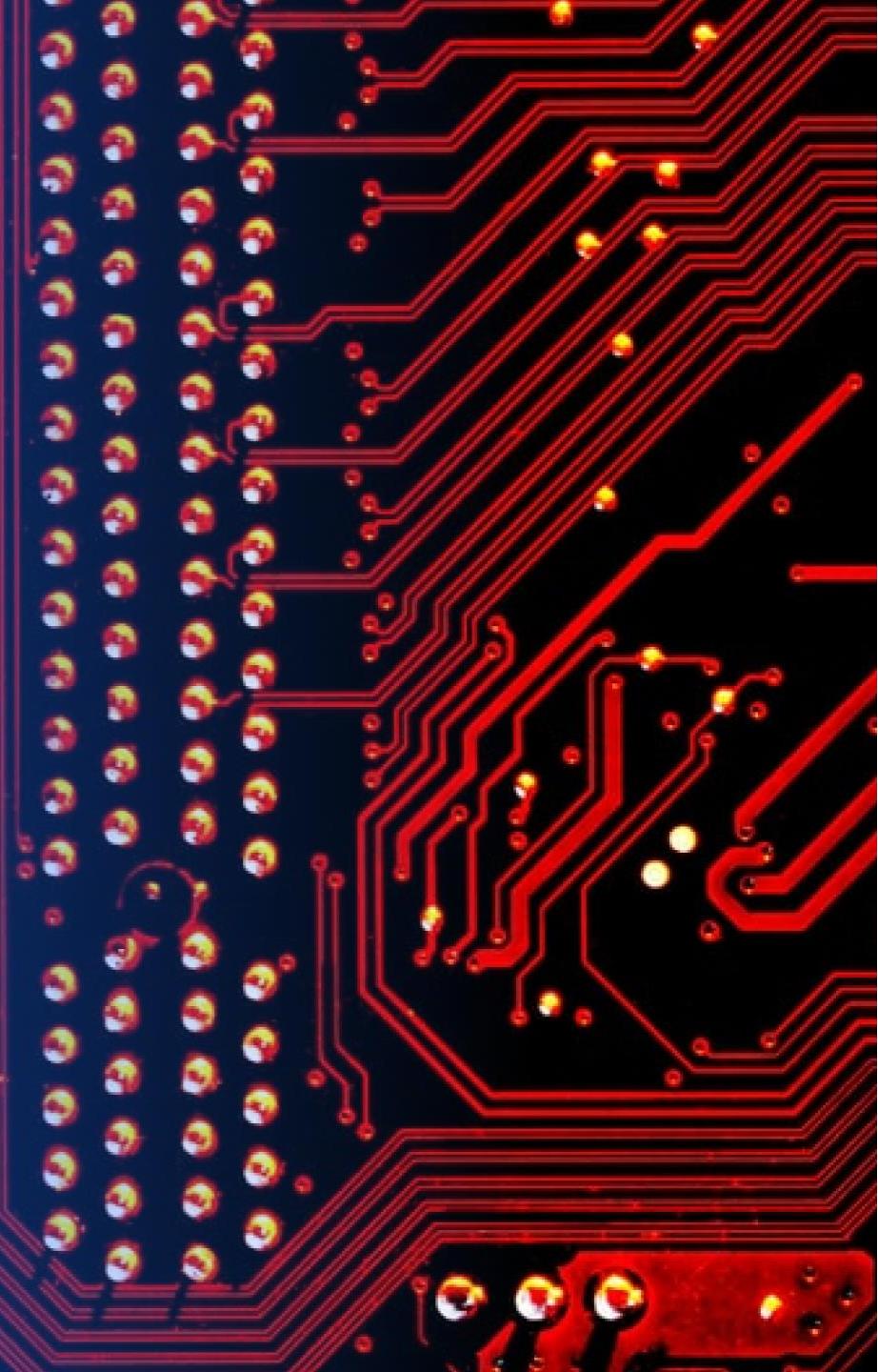


SITES PROXIMITY TO LANDMARKS



Section 4

Build a Dashboard with Plotly Dash

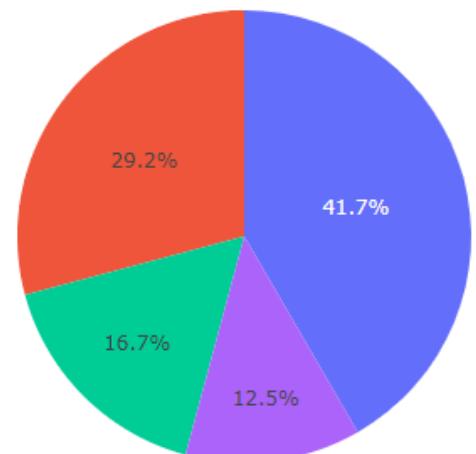


DASHBOARD – SUCCESS LAUNCHES FOR ALL SITES

SpaceX Launch Records Dashboard

All Sites

Total Success Launches By Site



- THE PIE CHART SHOWS THE SUCCESS RATE FOR THE 4 DIFFERENT SITES
- WE CAN SEE THAT KSC LC-39A HAS THE HIGHEST SUCCESS RATE

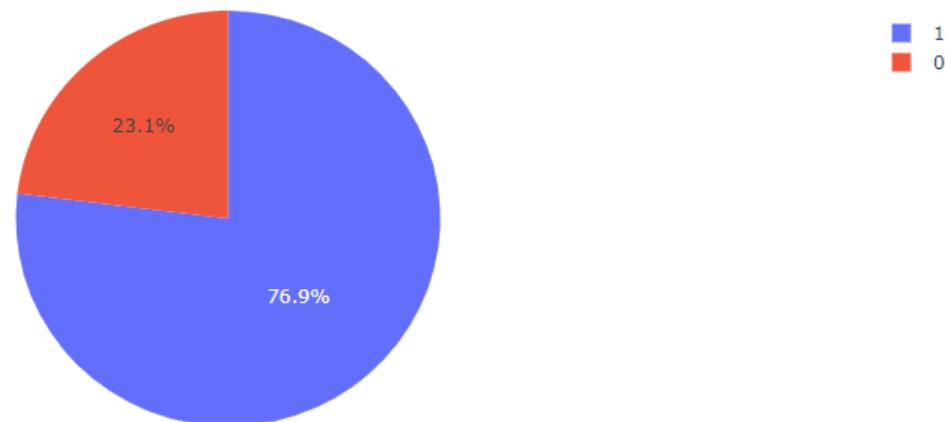
DASHBOARD – 'KSC LC-39A' SUCCESS RATIO

SpaceX Launch Records Dashboard

KSC LC-39A

x ▾

Total Success Launches for site KSC LC-39A



- THE SITE HAS ACHIEVED A 76.9% SUCCESS RATE

DASHBOARD – SCATTER PLOT OF PAYLOAD VS LAUNCH OUTCOME FOR ALL SITES



- LIGHT PAYLOAD FOR ALL SITES

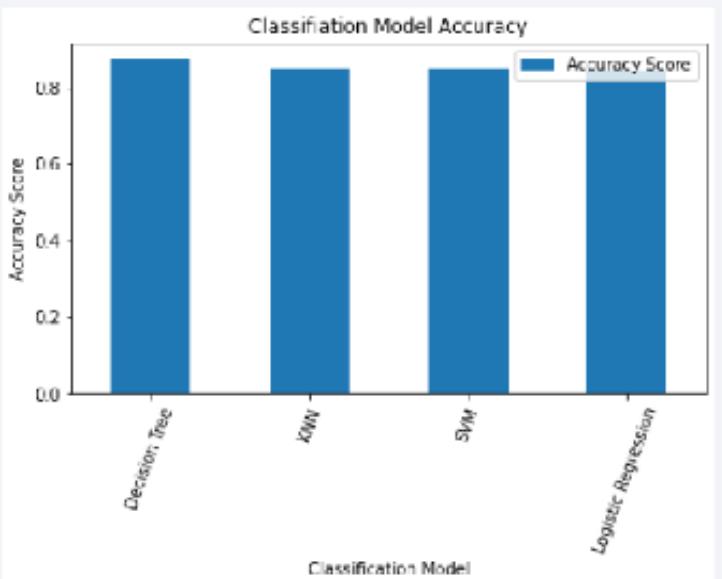
- HEAVY PAYLOAD FOR ALL SITES

The background of the slide features a dynamic, abstract design. It consists of several thick, curved lines that transition from a bright yellow at the top right to a deep blue at the bottom left. These lines create a sense of motion and depth, resembling a tunnel or a stylized landscape. The overall effect is modern and professional.

Section 5

Predictive Analysis (Classification)

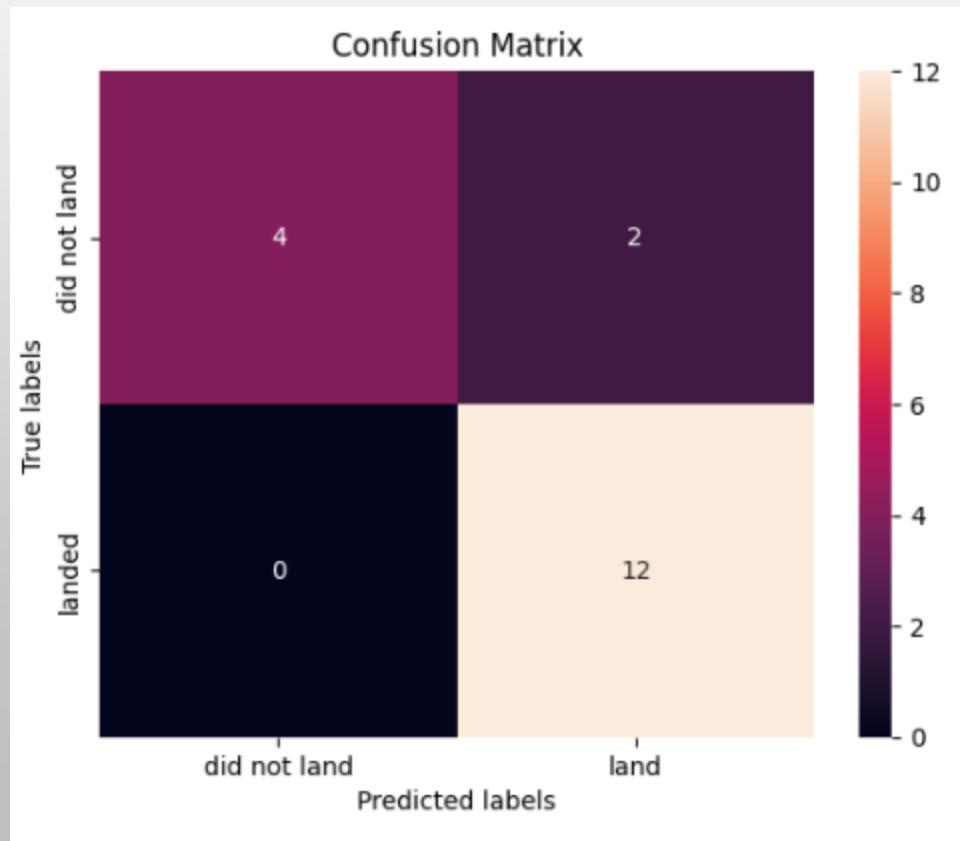
CLASSIFICATION ACCURACY



- BASED ON THE ACCURACY SCORE, AND AS SEEN ON THE BAR CHART, WE KNOW THE DECISION TREE IS THE MOST ACCURATE MODEL WITH A SCORE OF 0.8750

	Algo Type	Accuracy Score	Test Data Accuracy Score
2	Decision Tree	0.875000	0.833333
3	KNN	0.848214	0.833333
1	SVM	0.848214	0.833333
0	Logistic Regression	0.846429	0.833333

CONFUSION MATRIX



- THE CONFUSION MATRIX FOR THE DECISION TREE MODEL MADE A TOTAL OF **18** PREDICTIONS
- WE HAD **12** TRUE POSITIVES AND **4** TRUE NEGATIVES
- WE ALSO HAD **0** FALSE NEGATIVES AND **2** FALSE POSITIVES
- OVERALL, THE CLASSIFIER WAS CORRECT ABOUT **89%** OF THE TIME $((TP + TN)/TOTAL)$ AND HAD A MISCLASSIFICATION $((FN + FP)/TOTAL)$ OF **11%**

CONCLUSIONS

- AS THE NUMBER OF FLIGHTS INCREASE, SO DOES THE LIKELIHOOD OF A SUCCESSFUL FIRST STAGE LANDING
- SUCCESS RATE SEEMS TO GO UP AS PAYLOAD INCREASES, BUT THERE IS NO CONCLUSIVE CORRELATION BETWEEN PAYLOAD MASS AND SUCCESS RATE
- LAUNCH SUCCESS RATE INCREASED BY ABOUT 80% FROM 2013 TO 2020
- LAUNCH SITE 'KSC LC-39A' HAS THE HIGHEST LAUNCH SUCCESS RATE
- ES-L1, GEO, HEO AND SSO LEAD THE ORBIT TYPE SUCCESS RATE
- LAUNCH SITES ARE STRATEGICALLY LOCATED AWAY FROM CITIES AND CLOSE TO COASTLINES, RAILROADS AND HIGHWAYS, WHICH OBVIOUSLY MINIMIZES THE RISK OF CASUALTIES DUE TO A LANDING FAILURE
- THE BEST PERFORMING ML CLASSIFICATION MODEL IS THE DECISION TREE WITH AN ACCURACY OF ABOUT 87.5%. DESPITE THE RELATIVELY HIGH SCORE, MORE DATA MAY BE NEEDED IN ORDER TO BETTER FINE TUNE PREDICTIONS

APENDIX

Thank you!

