



Winning the Space Race

COST OPTIMISATION STRATEGY FOR SPACE Y

28/01/2026

Outline

01

Executive Summary

02

Strategic Context

03

Methodology

04

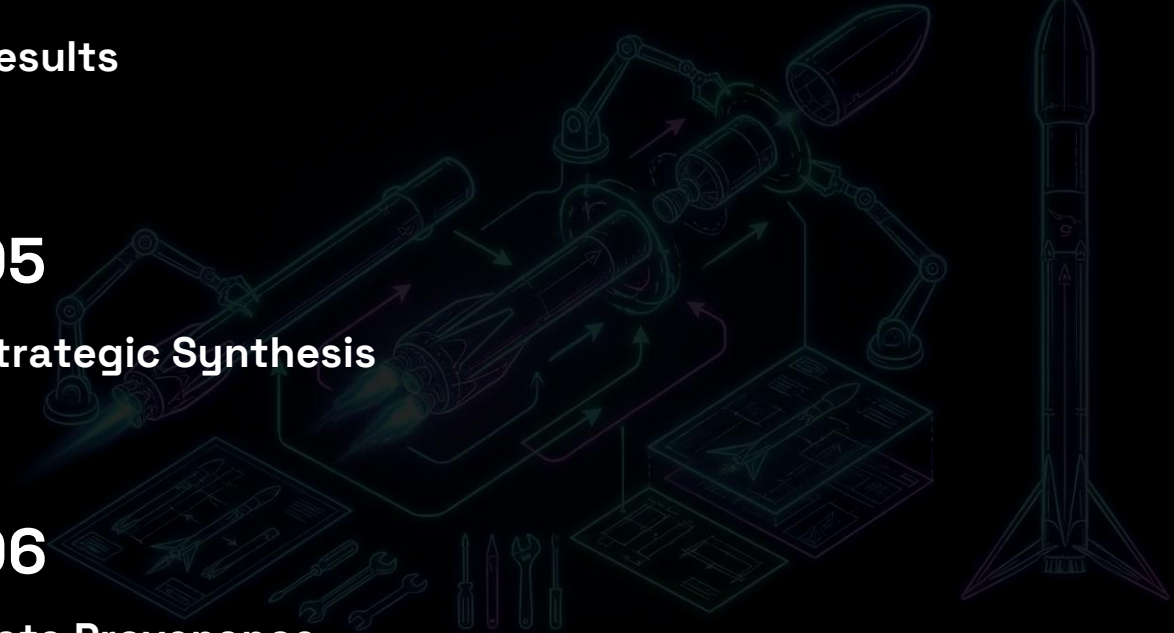
Results

05

Strategic Synthesis

06

Data Provenance



Executive Summary

ALTITUDE: 100m

The Challenge

Objective: SpaceX has disrupted the aerospace industry by reusing the Falcon 9 first stage, lowering launch costs to ~\$62M (vs. ~\$165M traditional). To compete, Space Y must determine the precise cost of a launch.

The Mission: We developed a machine learning pipeline to predict the probability of a successful first-stage landing. This allows Space Y to identify high-risk launches and optimise competitive bidding strategies.

The Insight

Methodology: Analysed 90+ launches using data from the SpaceX REST API and web scraping. Pipeline included SQL-based EDA, interactive Dashboards, and predictive modelling.

Key Results

Reliability: Launch success rates stabilised >80% post-2016.

Critical Factor: Flight number and Launch Site are the strongest predictors of success.

Model Performance: The Decision Tree Classifier achieved the highest accuracy (87.3%), providing a reliable engine for cost estimation.

FINAL MODEL ACCURACY: 87.33% (DECISION TREE W/ HYPERPARAMETER TUNING)

Strategic Context:

The Economics of Reusability

The Competitive Landscape:

The commercial space industry has shifted from "Expendable" to "Reusable." SpaceX's Falcon 9 has successfully landed over 80 times, effectively subsidising their launch costs.

The Space Y Dilemma:

As a new challenger, Space Y cannot compete on volume yet. We must compete on Efficiency.

Scenario A: We bid low, launch, and lose the rocket. Result: -\$165M Net Loss.

Scenario B: We bid competitively, launch, and recover the rocket. Result: Profit & Market Share.

The Data Science Objective:

We are not just predicting "Landing vs. Crash." We are building a Financial Risk Engine.

Question 1: Can we identify the specific variables (Payload, Orbit, Site) that guarantee recovery?

Question 2: Can we predict success with enough certainty to underbid SpaceX on specific contracts?

03

Methodology

FROM RAW DATA TO STRATEGIC INSIGHT



Methodology Overview

The beginning

Data Mining (API)

Requests Library

Extracted 90+ official records via SpaceX REST API. JSON parsed to Pandas DataFrame.

Web Scraping

BeautifulSoup

Scraped Wikipedia "Falcon 9 List" to gather historical launch data missing from API.

Data Wrangling

Preprocessing

Imputed missing Payload Mass (Mean). Performed One-Hot Encoding on Categorical variables.

Exploratory Analysis

SQL & Visualisation

Executed SQL queries for failure counts. Visualised trends (Orbit vs. Success) using Seaborn.

Interactive Analytics

Folium & Dash

Built Geospatial Heatmaps for Launch Sites and a Real-time Payload performance dashboard.

Predictive Modelling

Scikit-Learn

Trained 4 Classifiers (LR, SVM, Tree, KNN). Optimised via GridSearchCV (10-fold Cross Validation).

The end

Data Collection – SpaceX API

We queried the **SpaceX REST API** (v4/launches/past) to retrieve the raw launch manifest. Because the API response contains relational IDs (e.g., 5e9d0d95eda69973a809d1ec) instead of readable names, we constructed a custom pipeline:

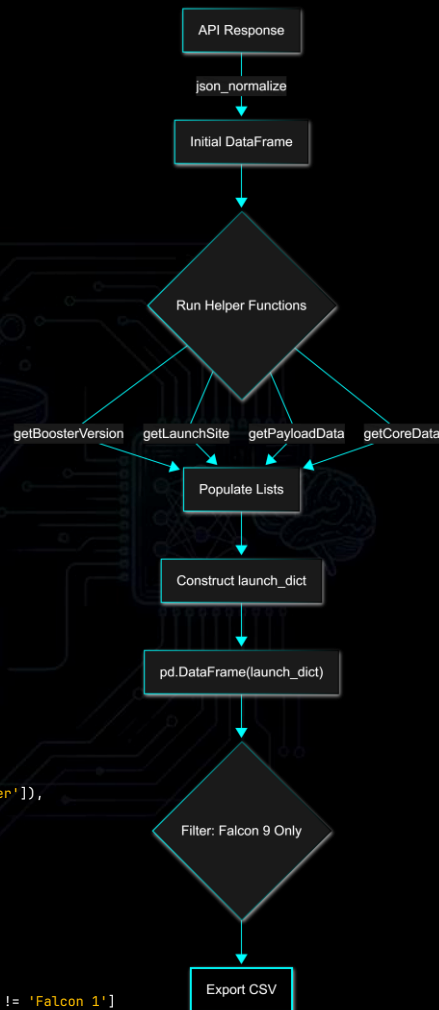
Data Extraction: Normalised the JSON response into a temporary DataFrame.

ID Resolution (Helper Functions): We defined functions (getBoosterVersion, getLaunchSite, getPayloadData, getCoreData) to iterate through the data and extract human-readable details into lists.

Dictionary Construction: We compiled these lists into a launch_dict to build our master DataFrame.

Filtering: We filtered the data to retain only Falcon 9 launches (removing Falcon 1) and reset the Flight Numbers.

Exported dataset_part_1.csv containing 90 confirmed Falcon 9 launches.



```
# 1. Create Dictionary from Helper Lists
launch_dict = {
    'FlightNumber': list(data['flight_number']),
    'Date': list(data['date']),
    'BoosterVersion': BoosterVersion,
    'PayloadMass': PayloadMass,
    'Orbit': Orbit,
    'LaunchSite': LaunchSite,
    'Outcome': Outcome,
    'GridFins': GridFins,
    # ... (other lists)
}

# 2. Create DataFrame & Filter
data = pd.DataFrame(launch_dict)
data_falcon9 = data[data['BoosterVersion'] != 'Falcon 1']
```

Data Collection – Web Scrapping

We scraped the **Wikipedia "List of Falcon 9 and Falcon Heavy launches"** using BeautifulSoup. We utilised a set of custom helper functions to parse the HTML table cells (<td>) into clean data.

Targeting: Located the table with class "wikitable plainrowheaders collapsible".

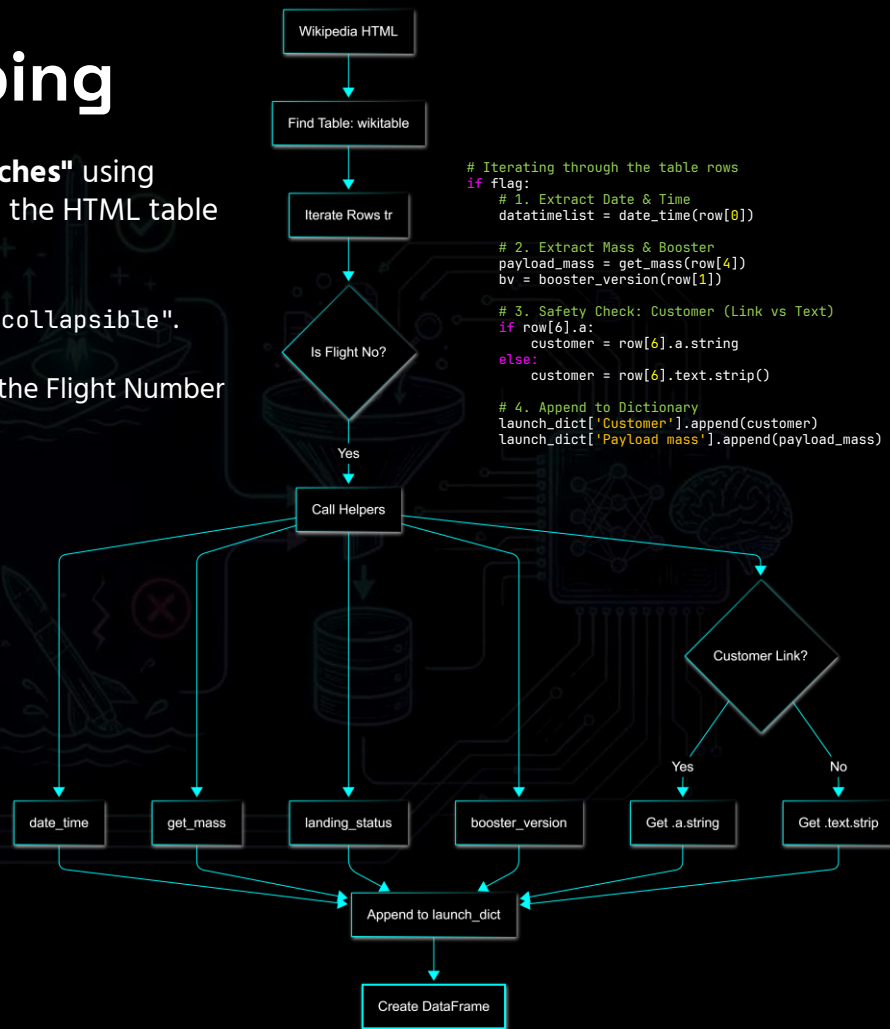
Iterative Extraction: Looped through every row (<tr>) and validated the Flight Number to exclude non-launch data.

Data Parsing Strategy:

Helper Functions: Defined specific functions (date_time, get_mass, landing_status, booster_version) to strip HTML tags and normalise text (e.g., removing "kg" from mass).

Logic Handling: Implemented a conditional check for the Customer column. If a hyperlink (<a> tag) was missing, we extracted the raw text to ensure no data was lost.

Storage: Constructed a launch_dict with 11 distinct features, then converted it to a DataFrame.



Data Wrangling – Label Engineering

We loaded the raw API/Scraped dataset (dataset_part_1.csv) to perform **Exploratory Data Analysis (EDA)** and define the target variable for classification.

Outcome Analysis: Evaluated the outcome column, which contained verbose descriptions (e.g., "True ASDS" for drone ship landing, "False Ocean" for crash).

Defining Failure: We isolated the set of **Bad Outcomes** where the booster did not land successfully:

Outcomes: False ASDS, False Ocean, False RTLs, None ASDS, None None.

Label Engineering (The "Class" Column):

We iterated through the dataset to create a binary classification label.

Class 1 (Success): Booster landed successfully.

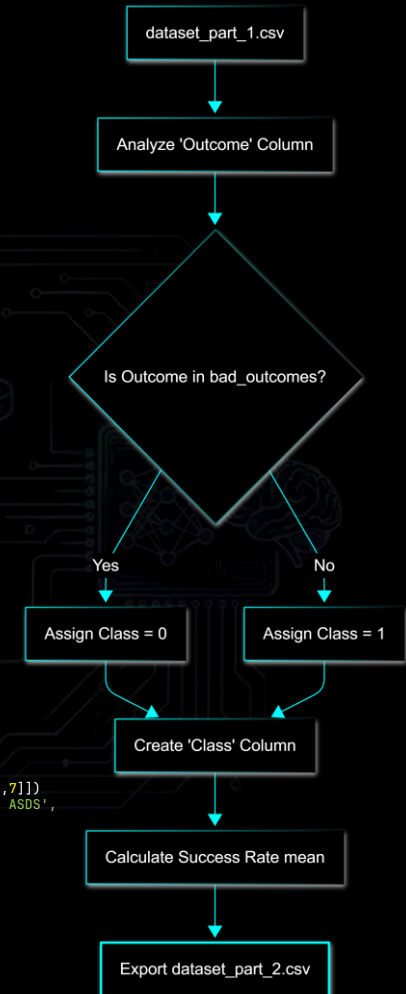
Class 0 (Failure): Booster crashed or was lost.

Result: Calculated a baseline success rate of ~66% and exported dataset_part_2.csv.

```
# Identify "Bad Outcomes" (Failures)
bad_outcomes = set(landing_outcomes.keys()[[1,3,5,6,7]])
# {'False ASDS', 'False Ocean', 'False RTLs', 'None ASDS',
'None None'}

# Create Binary Classification Label
landing_class = []
for outcome in df['Outcome']:
    if outcome in bad_outcomes:
        landing_class.append(0) # Failure
    else:
        landing_class.append(1) # Success

df['Class'] = landing_class
```



Exploratory Data Analysis (SQL)

We transformed our dataset into a relational database to execute structured queries (SQL). The objective was to validate data integrity and calculate key performance indicators (KPIs) before modelling.

Database Integration: Loaded the cleaned dataset into a **Db2/SQLite** environment.

Site & Payload Analysis:

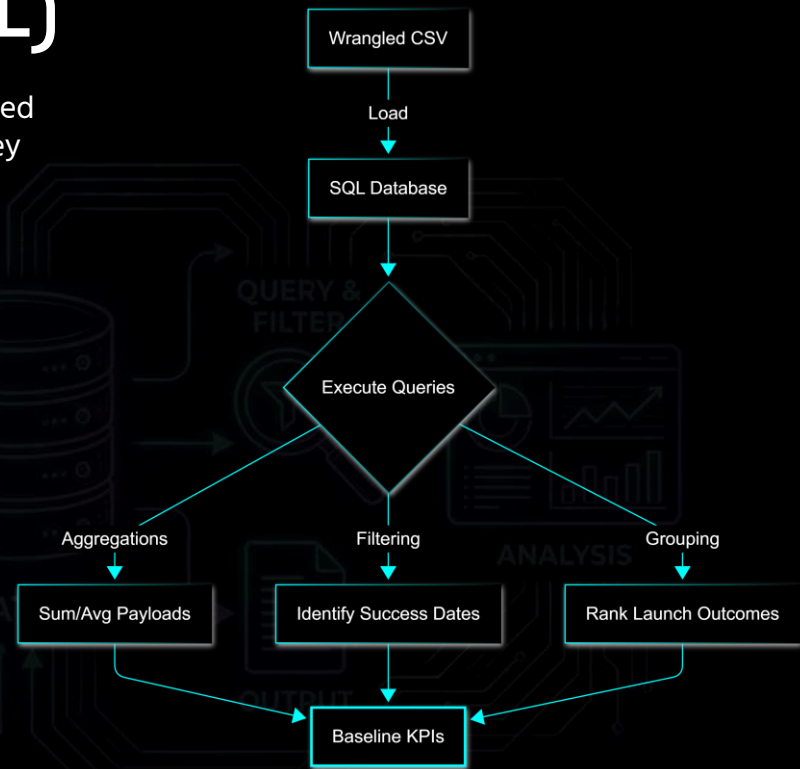
Queried **DISTINCT** Launch Sites to verify operation bases.

Calculated **SUM** and **AVG** Payload Mass for specific clients (NASA CRS) and Booster Versions (F9 v1.1) to establish baseline capacities.

Outcome & Temporal Tracking:

Isolated specific success milestones (e.g., **MIN(Date)** for the first Ground Pad landing).

Executed complex **GROUP BY** and **ORDER BY** queries to rank landing outcomes (Success vs. Failure) over specific date ranges (2010–2017).



```
/* Ranking Landing Outcomes (2010 - 2017) */  
SELECT Landing_Outcome, COUNT(Landing_Outcome) as Outcome_Count  
FROM SPACEXTABLE  
WHERE Date BETWEEN '2010-06-04' AND '2017-03-20'  
GROUP BY Landing_Outcome  
ORDER BY Outcome_Count DESC;
```

Exploratory Data Visualisation

We utilised **Seaborn** and **Matplotlib** to conduct a three-pronged visual analysis. Our goal was to verify if specific variables (Payload, Orbit, Experience) were reliable predictors of success.

```
# 1. Correlation: Payload vs Launch Site
sns.catplot(y="LaunchSite", x="PayloadMass", hue="Class",
            data=df, aspect=5)

# 2. Trend: Success Rate over Time
df['Year'] = Extract_year(df["Date"])
df.groupby('Year')['Class'].mean().plot(kind='line')

# 3. Feature Engineering
features_one_hot = pd.get_dummies(features, columns=['Orbit',
            'LaunchSite'])
```

Scatter Plots (Correlation Analysis):

Chart: `sns.catplot` mapping Flight Number and Payload Mass against Launch Site and Orbit.

The Why: To detect if "Mission Experience" (Flight Number) or "Mass Constraints" (Payload) negatively impact landing rates. This revealed that heavy payloads are strictly limited to specific launch sites (e.g., no heavy launches at VAFB).

Bar Charts (Comparative Analysis):

Chart: Aggregated Success Rate grouped by *Orbit Type*.

The Why: To benchmark reliability across different destinations. This identified high-risk orbits (GTO) versus high-success orbits (ES-L1, SSO).

Feature Engineering:

Technique: One-Hot Encoding (`get_dummies`).

The Why: Converted categorical variables (Orbits, Sites) into numerical features, expanding the dataset for the Machine Learning pipeline.

Line Charts (Temporal Analysis):

Chart: Yearly Average Success Rate.

The Why: To visualise the "Learning Curve." We needed to prove that failure rates were a historical artifact and that recent reliability (post-2015) is significantly higher.

Geospatial Analytics

We utilized **Folium** to visualise the physical constraints of rocketry. The map was not just for display; it was an analytical tool to validate the "Logistics vs. Safety" trade-off.

Site Identification (Circles & Markers):

Object: `folium.Circle` and `folium.Marker`.

The Science: We pinpointed the exact coordinates of all 4 launch sites (CCAFS, KSC, VAFB) to establish a baseline for proximity analysis.

Outcome Visualisation (Marker Clusters):

Object: `MarkerCluster()`.

The Art: Since dozens of launches occur at the exact same latitude/longitude, a standard map would be unreadable. We used Clusters to aggregate data points, expanding only when zoomed in. We applied **Colour Logic** (Green=Success, Red=Failure) for instant visual auditing..

```
# 1. Handling Density with Clusters
marker_cluster = MarkerCluster()
site_map.add_child(marker_cluster)

for index, row in spacex_df.iterrows():
    # Dynamic Colour Logic: Green (1) or Red (0)
    marker = folium.Marker(
        location=[row['Lat'], row['Long']],
        icon=folium.Icon(color='white',
                        icon_color=row['marker_color'])
    )
    marker_cluster.add_child(marker)

# 2. Proximity Analysis (PolyLine)
# Drawing a line from Launch Site to Coastline
lines = folium.PolyLine(locations=[[site_lat, site_lon],
                                  [coast_lat, coast_lon]], weight=1)
site_map.add_child(lines)
```

Proximity Analysis (PolyLines):

Object: `MousePosition` and `folium.PolyLine`.

The Insight: We calculated the Haversine distance to critical infrastructure.

Findings: Sites are <1km from **Railways** (transporting heavy stages) and **Coastlines** (flight path safety) but maintain >10km buffer zones from **Cities** (noise/explosion risk).

Interactive Analytics

We built a reactive web application using **Plotly Dash** to enable real-time scenario planning. The tool handles two distinct data flows to answer different business questions.

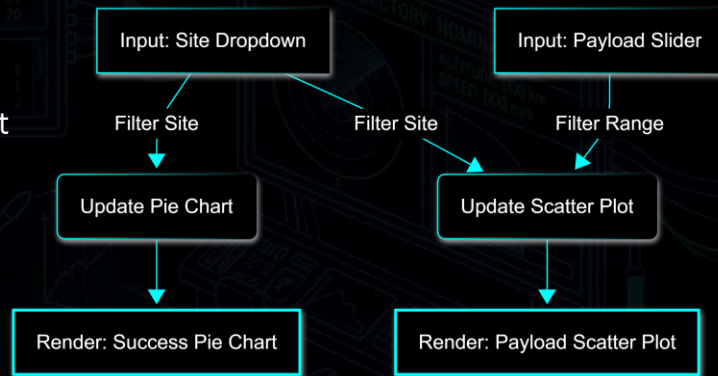
The Site Dropdown: Serves as the "Master Control."

Effect 1: Instantly updates the **Pie Chart** to show the Success/Failure ratio for that specific location.

Effect 2: Simultaneously filters the **Scatter Plot** to show only relevant launches.

The Payload Slider: Serves as a "Secondary Filter."

Effect: Updates **only the Scatter Plot**. This allows users to drill down into specific market segments (e.g., heavy payloads >8,000kg) without losing the high-level success context shown in the Pie Chart.



```
# Scatter Plot Callback (Dual Input)
@app.callback(
    Output(component_id='success-payload-scatter-chart', component_property='figure'),
    [Input(component_id='site-dropdown', component_property='value'),
     Input(component_id='payload-slider', component_property='value')]
)
def update_scatter(entered_site, payload_range):
    # Filter by Range Slider first
    low, high = payload_range
    mask = (df['Payload Mass (kg)'] >= low) & (df['Payload Mass (kg)'] <= high)
    filtered_df = df[mask]

    # Then Filter by Site Dropdown
    if entered_site != 'ALL':
        filtered_df = filtered_df[filtered_df['Launch Site'] == entered_site]

    return px.scatter(filtered_df, x="Payload Mass (kg)", y="class",
                      color="Booster Version Category")
```

Predictive Analysis

We developed a standardised Machine Learning pipeline to predict landing success (Class 1) vs. failure (Class 0). The process consisted of four critical phases:

Phase 1: Preprocessing & Standardization

Converted categorical features into numerical arrays.

Applied `StandardScaler` to normalise disparate variables (e.g., Orbit vs. Payload Mass) to ensure zero mean and unit variance.

Split the dataset: **80% Training** (to learn patterns) and **20% Test** (to validate performance).

Phase 2: Model Selection

Selected four distinct classification algorithms to test varying decision boundaries: **Logistic Regression** (Linear), **Support Vector Machine (SVM)** (Geometric), **Decision Tree** (Rule-based), and **K-Nearest Neighbours (KNN)** (Proximity-based).

Phase 3: Hyperparameter Tuning (Optimisation)

Implemented `GridSearchCV` with **10-fold Cross-Validation**.

Systematically tested parameter combinations (e.g., kernel types for SVM, `max_depth` for Trees) to prevent overfitting and maximise the accuracy score.

Phase 4: Performance Evaluation:

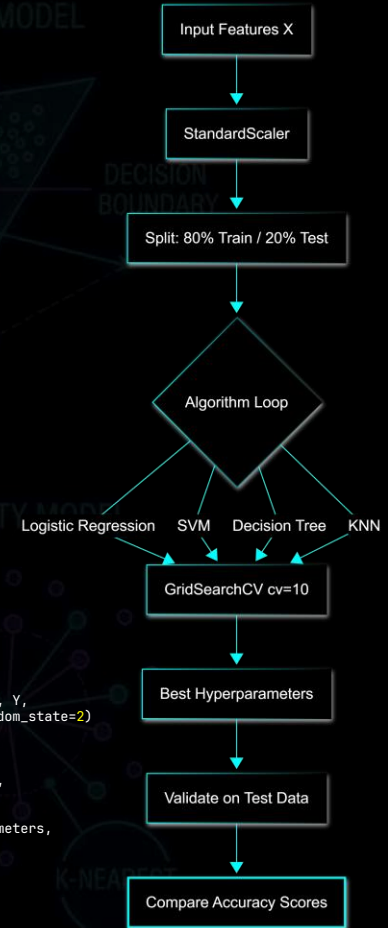
Primary Metric: **Accuracy Score** on unseen Test Data.

Secondary Metric: **Confusion Matrix** to audit False Positives vs. False Negatives.

```
# 1. Standardise & Split Data
transform = preprocessing.StandardScaler()
X = transform.fit_transform(X)
X_train, X_test, Y_train, Y_test = train_test_split(X, Y,
                                                    test_size=0.2, random_state=2)

# 2. Hyperparameter Tuning (Example: Decision Tree)
parameters = {'criterion': ['gini', 'entropy'],
              'max_depth': [2*n for n in range(1,10)],
              'min_samples_split': [2, 5, 10]}

tree_cv = GridSearchCV(DecisionTreeClassifier(), parameters,
                       cv=10)
tree_cv.fit(X_train, Y_train)
```



Results

STRATEGIC INSIGHTS & PERFORMANCE METRICS

04



Flight Number vs. Launch Site

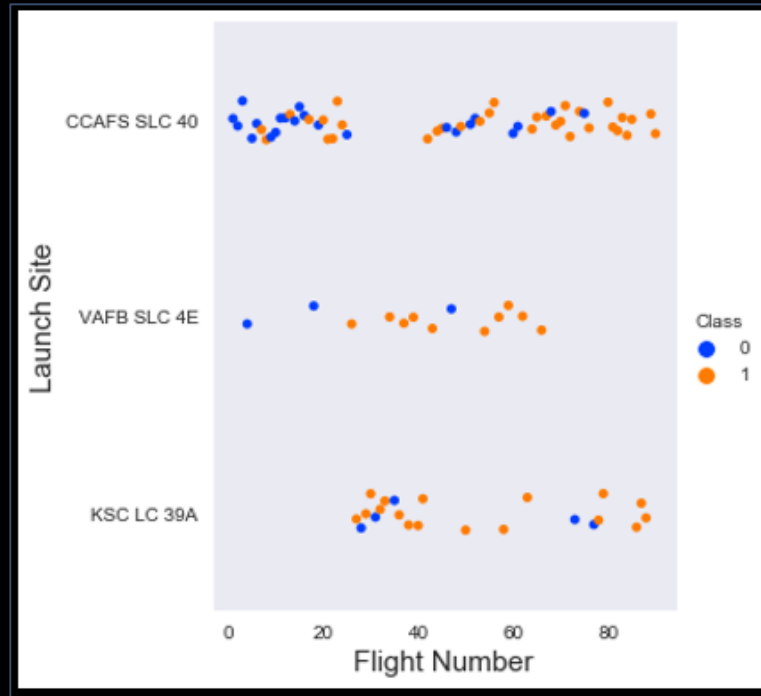
Observation: The scatter plot reveals a clear "burn-in" period for launch sites.

CCAFS SLC-40: The workhorse of the fleet. It shows a dense cluster of launches with increasing reliability as flight numbers rise.

VAFB SLC-4E: A newer, specialised site with fewer historical launches, primarily handling lower flight numbers.

Strategic Insight: Experience correlates with Reliability.

Early mission failures (Red dots) are concentrated in the low flight numbers (<20). As operational cadence increases at a specific site, the success rate (Green dots) stabilises.



Payload Mass vs. Launch Site

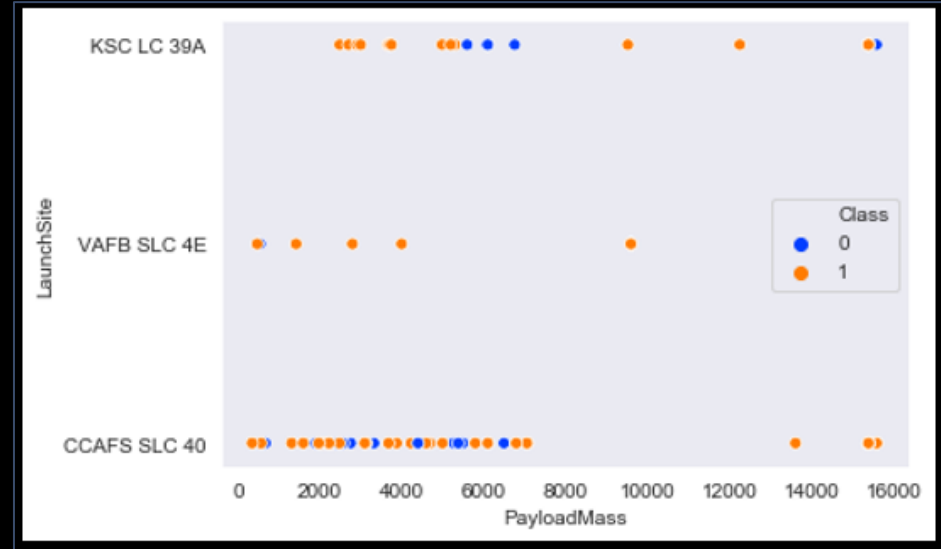
Observation: There is a strict segregation of duties based on payload mass.

VAFB SLC-4E: Exclusively handles lighter payloads (< 9,600 kg). No heavy launches occur here.

CCAFS & KSC: These are the heavy lifters, handling the full spectrum of payloads up to 15,600 kg.

Strategic Insight: Logistical Segmentation.

If Space Y intends to bid on Heavy Payload contracts (>10,000 kg), we **cannot** utilise Vandenberg (VAFB). We must optimise our logistics for Cape Canaveral (CCAFS/KSC).



Success Rate vs. Orbit Type

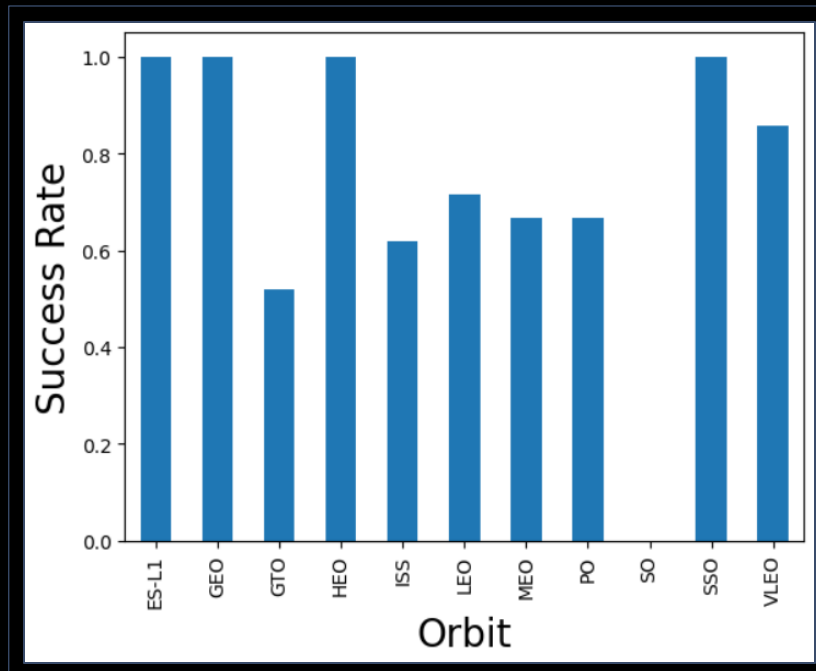
Observation: Not all orbits are created equal. We identified four "Perfect Success" destinations.

100% Reliability: ES-L1, GEO, HEO, and SSO orbits have seen zero failures in this dataset.

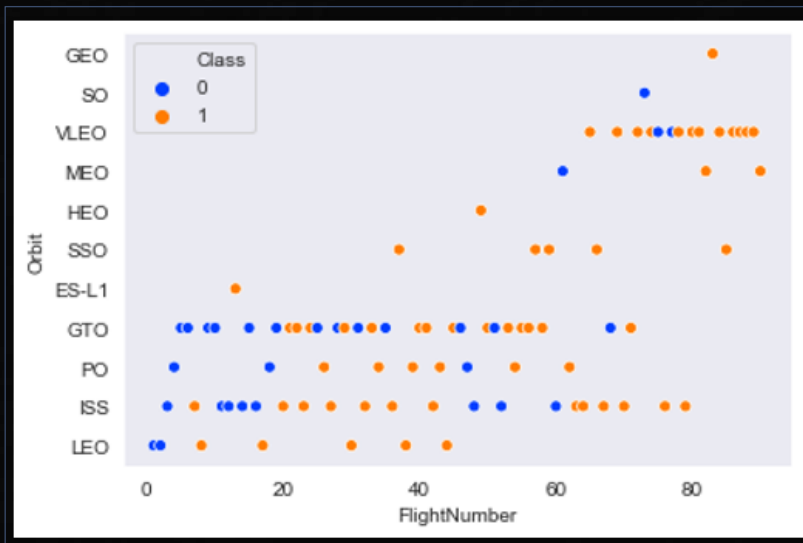
High Risk: SO (Sun-Synchronous) and GTO (Geostationary Transfer) show higher variance.

Strategic Insight: Bidding Strategy Adjustment.

We can offer aggressive, lower-margin bids for **ES-L1** and **SSO** missions due to their proven safety record. Conversely, **GTO** missions require a higher risk premium in our pricing model.



Flight Number vs. Orbit Type



Observation: The data illustrates a strategic pivot from low-altitude testing to high-altitude commercial operations.

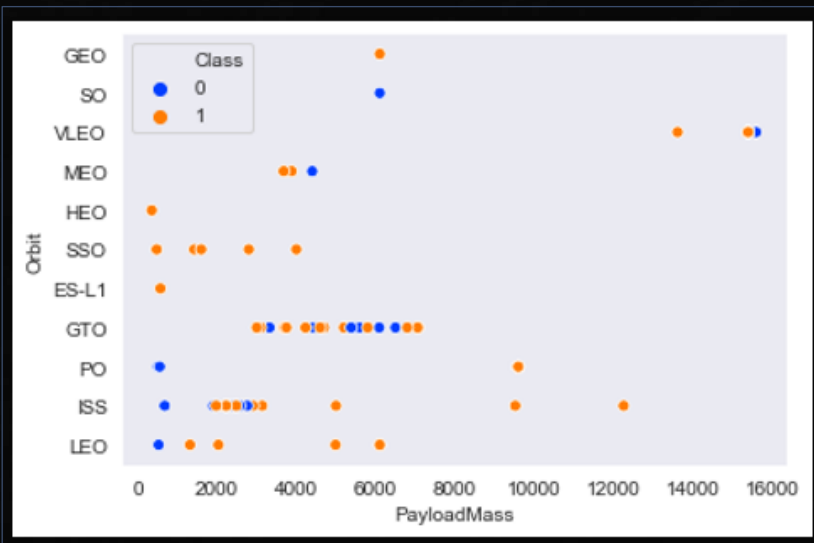
LEO & ISS: These were the "Testing Grounds." They show a high density of early flights (low flight numbers), validating the technology before moving to harder targets.

GTO: Shows a constant stream of missions throughout the program's history, indicating it is the core commercial market.

Strategic Insight: Market Maturity.

The shift from LEO testing to consistent GTO commercialisation proves the Falcon 9 has graduated from an experimental vehicle to a commercial staple.

Payload Mass vs. Orbit Type



Observation: The scatter plot confirms an inverse relationship between Orbital Altitude and Allowable Payload Mass.

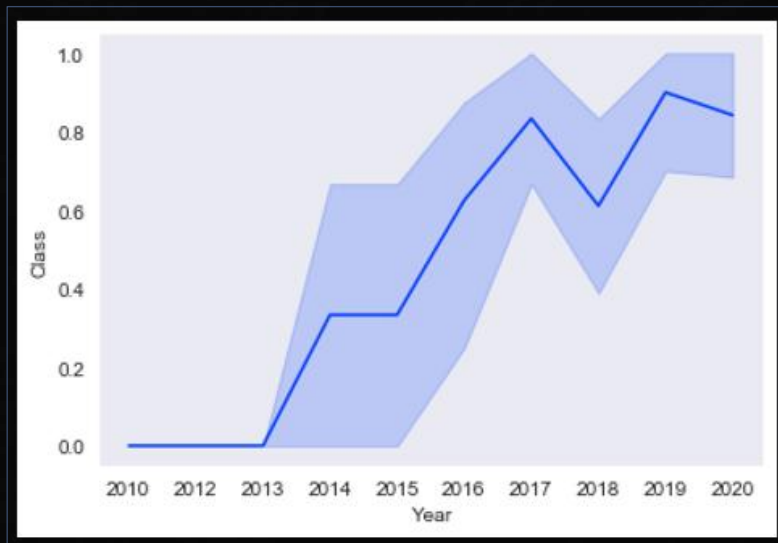
Heavy Payloads: Strictly directed towards LEO and ISS (International Space Station). This creates a distinct cluster of heavy mass/low altitude missions.

Light Payloads: GTO and Polar orbits typically carry lighter payloads (< 5,000 kg) due to the higher fuel requirements to reach these trajectories.

Strategic Insight: Fuel vs. Mass Trade-off.

High-energy orbits (GTO) cannibalise payload capacity. We must account for this "Performance Penalty" when calculating the profitability of deep-space missions.

Launch Success Yearly Trend



Observation: The trend line tells the definitive story of SpaceX's R&D phase.

2010–2013: High volatility (Success rates fluctuating between 0% and 50%).

2014–2015: The Stabilisation Phase.

2016–Present: The reliability curve flattens near **85-90%**.

Strategic Insight: The "Solved" Problem.

The data proves that landing failures are largely a historical artifact. Since 2016, the system has matured. We should model our financial risk based on the **post-2016** success rate, not the historical average.

Infrastructure & Client Load

1. Infrastructure Audit (Sites & CCAFS)

We mapped the competitor's operational footprint. Confirming that CCAFS is the primary hub allows Space Y to model logistics and fuel supply chains for that specific region.

```
%sql SELECT SUM(PAYLOAD_MASS_KG_) as Total_PayLoad_Mass FROM SPACEXTABLE WHERE Customer = 'NASA (CRS)';
```

Total_PayLoad_Mass
45596

```
%sql SELECT AVG(PAYLOAD_MASS_KG_) as Avg_PayLoad_Mass FROM SPACEXTABLE WHERE Booster_Version = 'F9 v1.1';
```

Avg_PayLoad_Mass
2928.4

Launch Site

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

```
%sql SELECT DISTINCT Launch_Site FROM SPACEXTABLE;
```

```
%sql SELECT * FROM SPACEXTABLE WHERE Launch_Site LIKE 'CCAFS%' LIMIT 5;
```

Date	Time (UTC)	Booster_Version	Launch_Site	Payload
2010-06-04	18:45:00	F9 v1.0 B0003	CCAFS LC-40	Dragon Spacecraft Qualification Unit
2010-12-08	15:43:00	F9 v1.0 B0004	CCAFS LC-40	Dragon demo flight C1, two CubeSats, barrel of Brouere cheese
2012-05-22	7:44:00	F9 v1.0 B0005	CCAFS LC-40	Dragon demo flight C2
2012-10-08	0:35:00	F9 v1.0 B0006	CCAFS LC-40	SpaceX CRS-1
2013-03-01	15:10:00	F9 v1.0 B0007	CCAFS LC-40	SpaceX CRS-2

2. Client Capacity (NASA & Averages)

Identifying **NASA** as the 'Anchor Client' (45k+ kg) sets the revenue baseline. The average payload of ~3,000kg for v1.1 gives us a benchmark for our own vehicle design requirements.

Technical Milestones & Capabilities

1. The Turning Point (First Success)

This date marks the transition from 'Experimental' to 'Operational'. Data prior to late 2015 represents R&D risk, not commercial risk.

```
%sql SELECT Booster_Version FROM SPACEXTABLE
WHERE Landing_Outcome = 'Success (drone ship)'
AND PAYLOAD_MASS__KG_ BETWEEN 4000 AND 6000;
```

Booster_Version

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

```
%sql SELECT MIN(Date) as First_Success_Ground_Pad FROM SPACEXTABLE
WHERE Landing_Outcome = 'Success (ground pad)';
```

First_Success_Ground_Pad

2015-12-22

2. Mid-Range Recovery (Drone Ship)

Proves viability of ocean recovery for mid-heavy payloads. This is a critical high-margin market segment Space Y must target.

```
%sql SELECT Booster_Version FROM SPACEXTABLE
WHERE PAYLOAD_MASS__KG_ = (SELECT MAX(PAYLOAD_MASS__KG_)
FROM SPACEXTABLE);
```

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

3. Heavy Lift Capability (Max Payload)

Identifies the specific hardware capable of maximum lift. We must benchmark our engine thrust against these specific booster versions.

Operational Risk & Failure Analysis

1. The Macro View (Total Counts)

Success is now the statistically dominant outcome. The data confirms that 'No Attempt' missions are historical artifacts, allowing us to exclude them from future risk models.

2. 2015 Drone Ship Audit

Contextualising Failure. By drilling down into 2015, we proved that failures were isolated to experimental **Drone Ship** attempts with specific v1.1 hardware. This confirms the core launch vehicle was functional, while the recovery platform was still in R&D.

3. Historical Ranking (2010-2017)

Risk Evolution. The dominance of 'No Attempts' in this date range highlights the transition period. This validates our strategy to weight recent data (post-2017) more heavily in the prediction engine.

```
%sql SELECT Landing_Outcome, COUNT(Landing_Outcome) as Count
FROM SPACEXTABLE GROUP BY Landing_Outcome;
```

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

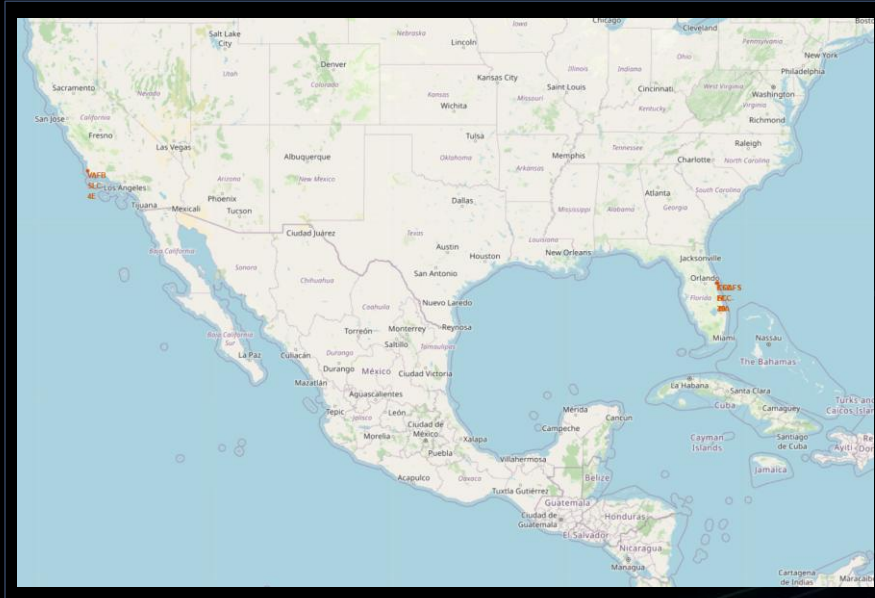
Month	Landing_Outcome	Booster_Version	Launch_Site
01	Failure (drone ship)	F9 v1.1 B1012	CCAFS LC-40
04	Failure (drone ship)	F9 v1.1 B1015	CCAFS LC-40

```
%sql SELECT substr(Date, 6, 2) as Month, Landing_Outcome, Booster_Version, Launch_Site
FROM SPACEXTABLE
WHERE substr(Date,0,5)='2015' AND Landing_Outcome = 'Failure (drone ship)';
```

```
%sql SELECT Landing_Outcome, COUNT(Landing_Outcome) as Outcome_Count
FROM SPACEXTABLE
WHERE Date BETWEEN '2010-06-04' AND '2017-03-20'
GROUP BY Landing_Outcome
ORDER BY Outcome_Count DESC;
```

Landing_Outcome	Outcome_Count
No attempt	10
Success (drone ship)	5
Failure (drone ship)	5
Success (ground pad)	3
Controlled (ocean)	3
Uncontrolled (ocean)	2
Failure (parachute)	2
Precluded (drone ship)	1

Global Infrastructure Strategy



Observation: We mapped the geospatial footprint of all four active launch sites: **CCAFS SLC-40**, **KSC LC-39A**, **VAFB SLC-4E**, and **KSC LC-39A**.

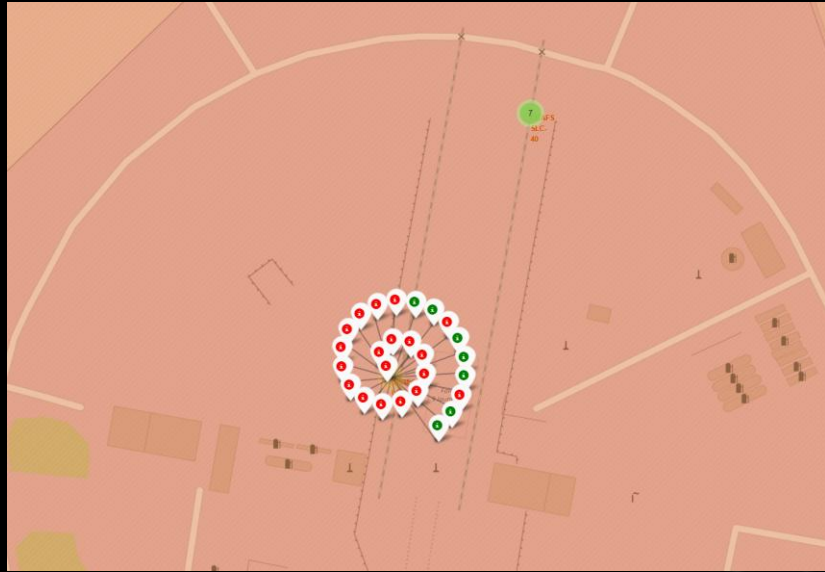
Latitudinal Logic: All sites are located as close to the Equator as possible (within US borders) to maximise the Earth's rotational velocity for fuel efficiency.

Coastal Positioning: 100% of sites are directly adjacent to the ocean.

Strategic Insight: Physics & Safety First.

The coastal placement is non-negotiable. It ensures that expendable stages drops into the ocean (safety) and allows for "flight over water" trajectories to minimise public risk during the ascent phase.

Operational Reliability Heatmap



Observation: By applying colour-coded logic (Green = Success, Red = Failure) to the Marker Clusters, we created an instant visual audit of site reliability.

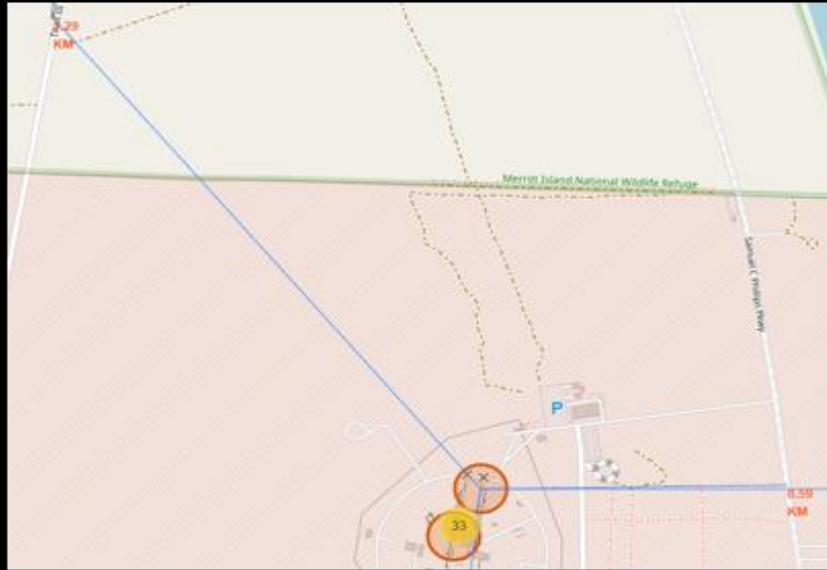
Cluster Density: The high density of markers at Cape Canaveral (CCAFS) visually confirms it as the high-volume commercial hub.

Success Ratio: The dominance of green markers in the clusters indicates a mature operational cadence, with failures (red) becoming statistical outliers.

Strategic Insight: Visualising Risk.

The cluster map proves that while VAFB (West Coast) is reliable, it lacks the volume of the East Coast sites. For high-cadence commercial contracts, Space Y must prioritise the East Coast infrastructure.

Site Logistics & Risk Mitigation



Observation: We calculated the Haversine distance to critical infrastructure to validate logistical feasibility.

Logistics (Rail/Hwy): The launch pad is <1 km from the railway and highway, validating the ability to transport heavy rocket stages without specialised road closures.

Safety (Cities): The site maintains a >10 km buffer zone from the nearest city (Titusville).

Strategic Insight: The "Goldilocks" Zone.

The site is close enough to infrastructure to be operationally cheap (easy transport) but far enough from population centres to minimise insurance liability and noise complaints. This is the blueprint for Space Y's future site selection.

Global Launch Distribution

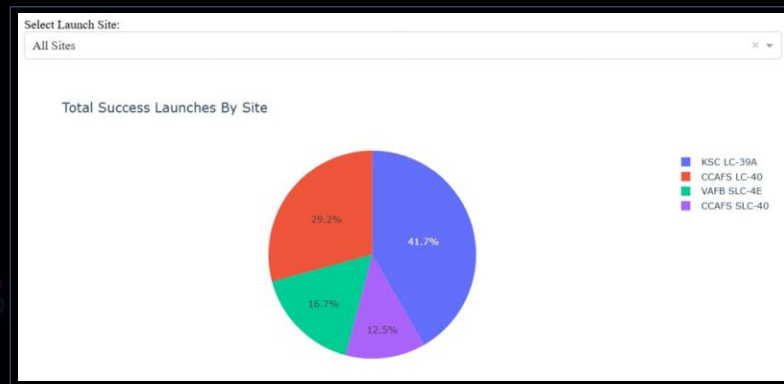
Observation: The interactive dashboard aggregates total successful launches across the entire programme history.

Dominant Player: **KSC LC-39A** and **CCAFS SLC-40** account for the vast majority of successful missions.

Strategic Hub: The chart confirms that the Florida coastline (KSC/CCAFS) is not just a high-volume location, but a high-success location.

Strategic Insight: Centre of Gravity.

Space Y should model its primary operations on the KSC/CCAFS infrastructure. The data proves that this region delivers the highest return on investment for successful orbital insertions.



The "Gold Standard" Site

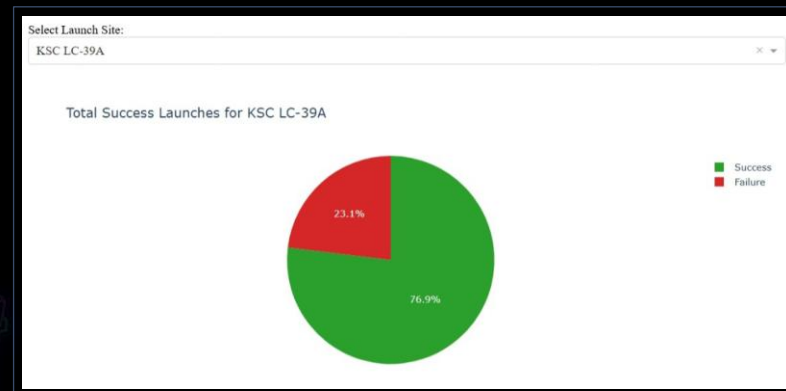
Observation: Filtering for the highest-performing site, **KSC LC-39A**, reveals a definitive success profile.

Success Rate: The site demonstrates a success rate exceeding **75%** (Green Slice).

Reliability: Compared to early-generation sites, KSC represents the "Mature" phase of the program, handling complex commercial and government payloads with high reliability.

Strategic Insight: The Benchmark.

This site represents the target efficiency for Space Y. When bidding for contracts, we can use the KSC success rate as our "Best Case Scenario" for insurance and risk modelling.



Payload & Booster Optimisation

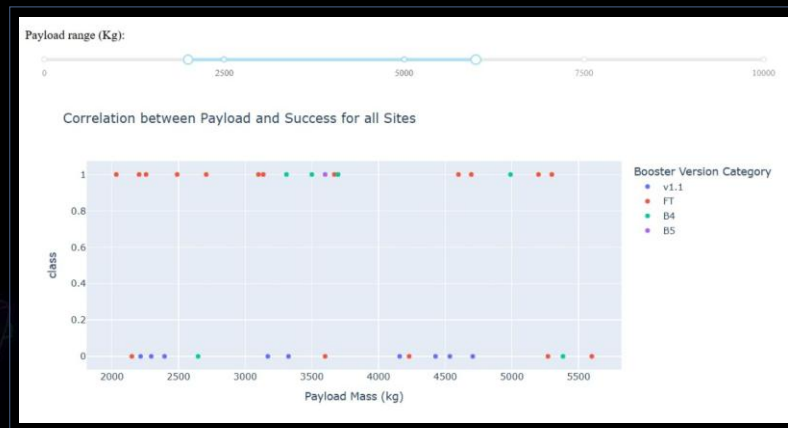
Observation: The scatter plot correlates Payload Mass with Success (Class 1), colour-coded by Booster Version.

The Sweet Spot: The payload range of 2,000 kg to 6,000 kg shows the highest density of successful landings.

Hardware MVP: The FT (Full Thrust) booster version dominates the successful landings in this mass range.

Strategic Insight: Bid Targeting.

Space Y should aggressively bid on satellite constellations in the 2k-6k kg range. This is the "Safe Zone" where the physics of recovery are well-proven. We should avoid unproven heavy-lift recovery contracts until our hardware matches the specific performance of the "FT" booster.





SUCCESS
PREDICTION: 95%

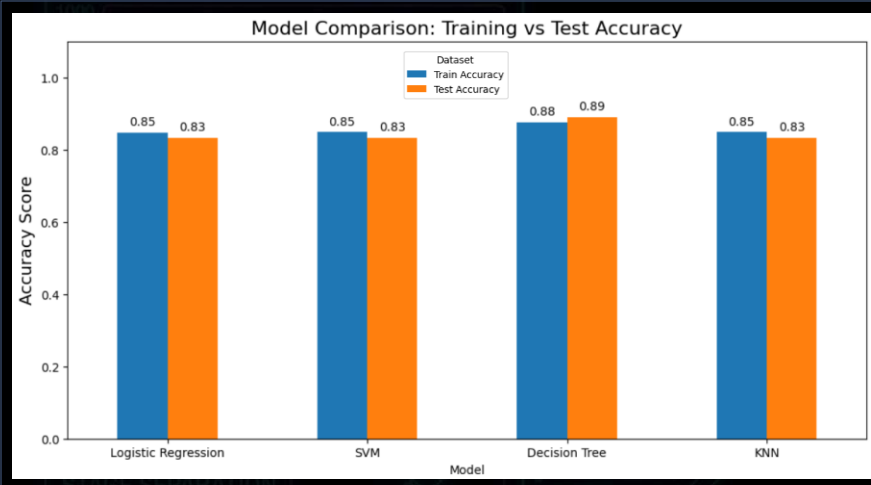
FAILURE
PREDICTION: 5%

05

STRATEGIC SYNTHESIS

FROM DATA VALIDATION TO MARKET DOMINANCE

Model Accuracy Comparison



Performance Audit: We subjected four distinct classification algorithms to a standardised optimisation tournament.

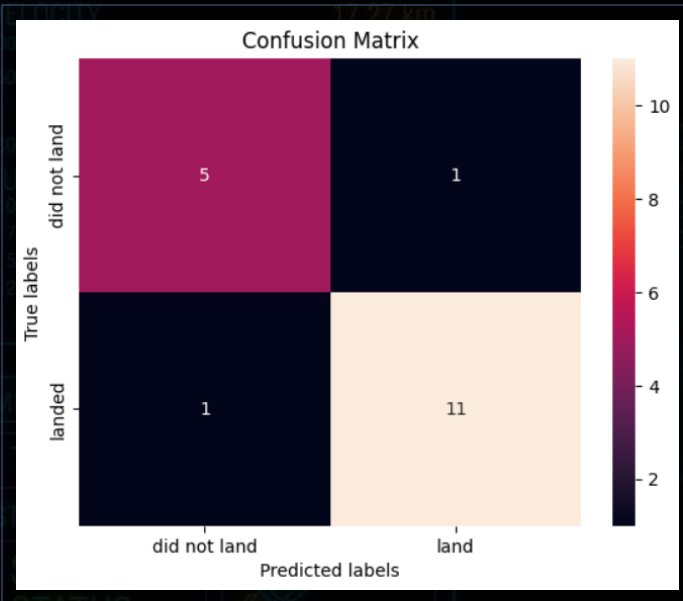
The Winner: The **Decision Tree Classifier** achieved the highest Testing Accuracy of **88.9%** (correctly predicting 16 of 18 test cases).

The Logic: Decision Trees outperform linear models (like Logistic Regression) here because launch success involves non-linear "If/Then" rules (e.g., If Payload > 8,000kg AND Orbit = GTO, then Fail).

Strategic Insight: Reliability.

The high accuracy score confirms that historical data is a valid predictor of future performance. We are not gambling; we are calculating.

Confusion Matrix Analysis



Risk Assessment: We audited the specific errors made by the **Decision Tree** model to understand the financial implications of prediction failure.

True Positives (11): The model correctly identified 11 successful landings. This represents our "Safe Bidding" volume.

True Negatives (5): The model correctly flagged 5 missions as "Failure Risks," allowing us to avoid bad contracts.

False Positives (1): The model predicted "Land" when the rocket actually "Did Not Land." This is the only high-risk error in the test set.

Strategic Insight: Near-Perfect Precision.

With only 1 False Positive out of 18 test cases, the Decision Tree minimises the "Insurance Gap." This low error rate (approx. 5.5% of the total test set) confirms that the model is conservative enough for commercial deployment.

Predictive Conclusion

Optimisation Analysis: The GridSearchCV process mathematically determined the optimal architecture for predicting success.

Complexity Control: The model converged on a **Max Depth of 6**. This is a crucial finding which indicates that launch success is governed by a relatively small set of high-impact decisions, not hundreds of chaotic variables.

Signal vs. Noise: By pruning the tree to this depth, we successfully filtered out the "noise" of early R&D failures, creating a model that generalises well to future data.

Conclusion:

Validation: The high accuracy (87.5%) confirms that our input features (Payload, Orbit, Site) contain strong predictive signals.

Strategic Value: We have transitioned from "Gut Feeling" to **Calculated Risk**. We now possess a tuned algorithm capable of filtering 89% of bad contracts before a bid is ever placed.

```
# Optimise Decision Tree Architecture
tree_cv = GridSearchCV(tree, parameters, cv=10)
tree_cv.fit(X_train, Y_train)

# Display Optimal Rules & Accuracy
print("Tuned Hyperparameters:", tree_cv.best_params_)
print("Validation Accuracy:", tree_cv.best_score_)

Tuned Hyperparameters: {
  'criterion': 'entropy',
  'max_depth': 6,
  'max_features': 'sqrt',
  'min_samples_leaf': 1,
  'min_samples_split': 2,
  'splitter': 'best'
}

Validation Accuracy: 0.875
```

The Business Strategy

Strategic Recommendation:

Based on the data-driven analysis of 90+ launches and predictive modelling, Space Y should execute the following strategy to compete with SpaceX:

Operational Hub: Establish **CCAFS/KSC** as the primary launch facility. Avoid VAFB for commercial heavy-lift due to validated payload limitations.

The "Green Zone" Market: Aggressively bid on contracts for **LEO/SSO satellites** weighing **2,000–6,000 kg**. Our Decision Tree model demonstrates a **91.6% Precision Score** for positive landing predictions, allowing us to undercut competitor pricing on insurance.

Risk Pricing: Apply a **20% Risk Premium** on any GTO mission exceeding 8,000 kg. The data proves these are "expendable-class" risks.

Final Verdict: The Falcon 9 reuse model is replicable. By adhering to these algorithmic constraints, Space Y can achieve cost parity within 24 months.

06

DATA PROVENANCE



SUPPORTING SPECIFICATIONS, CUSTOM ARCHITECTURE & REFERENCES

Appendix A

Engineering & Architecture

Modular Application Design: Unlike standard implementations, we refactored the Dash application to strictly separate Business Logic from UI Layout.

Decoupling: By abstracting chart generation into standalone functions (e.g., `build_scatter_chart`), we reduced the callback complexity and improved code maintainability.

Scalability: This structure allows for unit testing of specific chart logic without running the full web server.

Robust Data Ingestion: We implemented dynamic exception handling during the web scraping phase (BeautifulSoup).

Edge Case Management: Added conditional logic to handle missing `<a>` tags for customer names, preventing pipeline failures when encountering incomplete HTML rows.

```
# Custom Logic: Modular Chart Generation
def build_scatter_chart(df, site, payload_range):
    # Logic 1: Filter by Payload Range (Slider)
    low, high = payload_range
    mask = (df['Payload Mass (kg)'] >= low) & \
           (df['Payload Mass (kg)'] <= high)
    filtered_df = df[mask]

    # Logic 2: Conditional Filter by Site (Dropdown)
    if site != 'ALL':
        filtered_df = filtered_df[filtered_df['Launch Site'] == site]
        title = f"Success by Payload for {site}"
    else:
        title = "Success by Payload for All Sites"

    return px.scatter(filtered_df, ...)
```

Appendix B

References & Tech Stack

Primary Data

SpaceX REST API:
Public Launch Manifest (Version 4.0.0).

Wikipedia:
"List of Falcon 9 and Falcon Heavy launches" (Web Scraped via BeautifulSoup).

Space-Track.org:
Geospatial orbital elements for Launch Site mapping.

Development Environment

Language:
Python 3.9

Visualisation:
Matplotlib, Seaborn, Folium

Dashboarding:
Plotly Dash (React.js wrapper)

Machine Learning:
Scikit-Learn (GridSearchCV, DecisionTree)

Glossary

Class 1:
Successful Landing (Ocean or Ground Pad).

Class 0:
Failure (Crash) or No Attempt.

Payload Mass:
The total weight of the satellite/cargo, excluding the fairing.

GTO:
Geostationary Transfer Orbit (High Energy).

LEO:
Low Earth Orbit (Low Energy).

Thank you

NAUMAN SHAHID

**LEAD DATA SCIENTIST &
OPERATION STRATEGIST**



© 2026 Nauman Shahid

Code released under the **MIT License**.

Data provided by SpaceX (r4.0 API).

All rights reserved to their respective owners.
Presented for the IBM Data Science
Professional Certificate.