# DATA ANALYSIS PORTFOLIO

By Mohan S

# PROFESSIONAL BACKGROUND

I am a committed Data Analyst with a solid background in data analysis, statistical modeling, and business intelligence tools. Holding a Master's in Computer Applications from The National Institute of Engineering, Mysore, and a Bachelor's in Computer Applications, I have acquired skills in SQL, Power BI, Tableau, Excel, and Python to analyze and interpret intricate datasets for data-driven decision-making.

Currently, I am working as a Data Analytics Trainee in Trainity, where I undertook several projects that were enhancing decision-making processes by pattern and trend identification, improvement in data quality through cleaning and organizing, and interactive dashboards that allowed for better data interpretation by 30%. These insights had been converted into actionable recommendations that have improved efficiency and stakeholder engagement.

Previously, I interned at Webblitz Softwares as a Python Developer, where I worked with SQL, Django, and web technologies, improving system performance and user experience. Additionally, I have worked on diverse data-driven projects such as:

- Hiring Process Analytics – Performed statistical analysis on hiring data, developed insights to enhance gender diversity hiring strategies, and created data visualizations to optimize recruitment efficiency.
- Instagram User Analytics – Conducted marketing analytics using SQL on 50,000+ user data points, identified engagement trends, and improved campaign targeting by 20%.

Beyond my technical expertise, I am skilled in business analysis, problem-solving, and stakeholder communication, ensuring that insights are both accurate and actionable. I am passionate about leveraging data analytics to drive business impact and continuously enhancing my analytical skills to solve real-world challenges.

# TABLE OF CONTENT

# TABLE OF CONTENT

# *INSTAGRAM USER ANALYTICS*

## Description

The goal of this project is to use SQL to analyze Instagram user data stored in a relational database. The insights derived aim to assist Instagram's product, marketing, and development teams in making informed decisions about user engagement, marketing strategies, and operational improvements. The analysis covers various areas such as identifying loyal users, inactive users, contest winners, popular hashtags, and user registration trends.

## Problem

The Instagram team seeks insights on user engagement, retention, and revenue growth. Your task is to identify top-performing content types, peak activity hours, and user retention trends. Additionally, analyze high-value user segments driving ad revenue and measure the impact of new feature releases. These insights will guide product development, marketing strategies, and monetization efforts.

## Design

- Imported & Merged Data – Combined datasets for analysis.
- Cleaned Data – Removed duplicates, blank cells, and incorrect values.
- Formatted Columns – Improved headers for better readability.
- Analyzed Data – Used SQL queries to extract insights on user engagement.
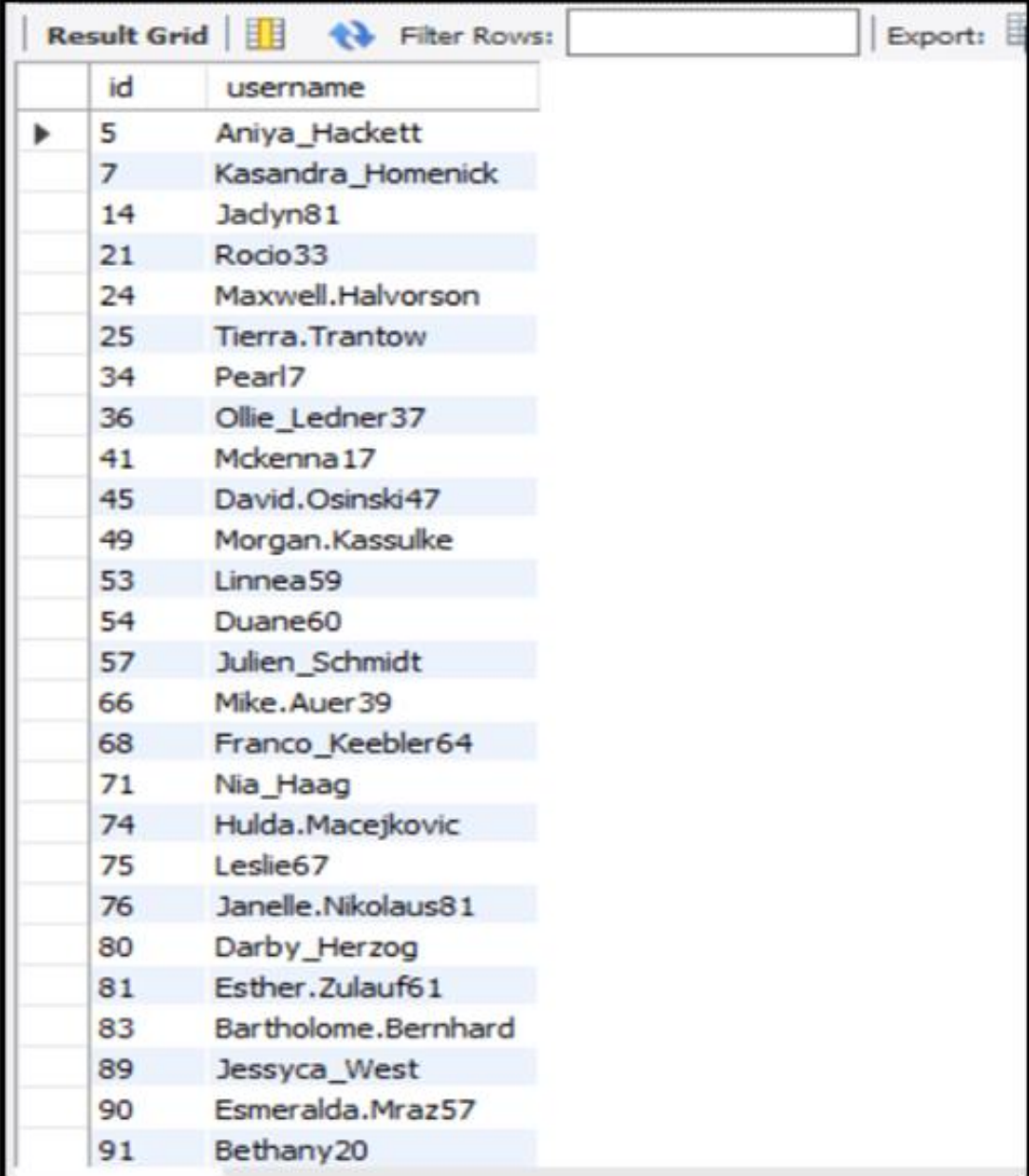- Provided Insights – Helped teams make data-driven decisions.

## Finding – 1

**Identify the five longest-active users on Instagram from the provided database for the marketing team's Loyal User Reward program.**



These users have been on Instagram the longest.

# Finding – 2

Identify users who have never posted on Instagram for the Inactive User Engagement email campaign.



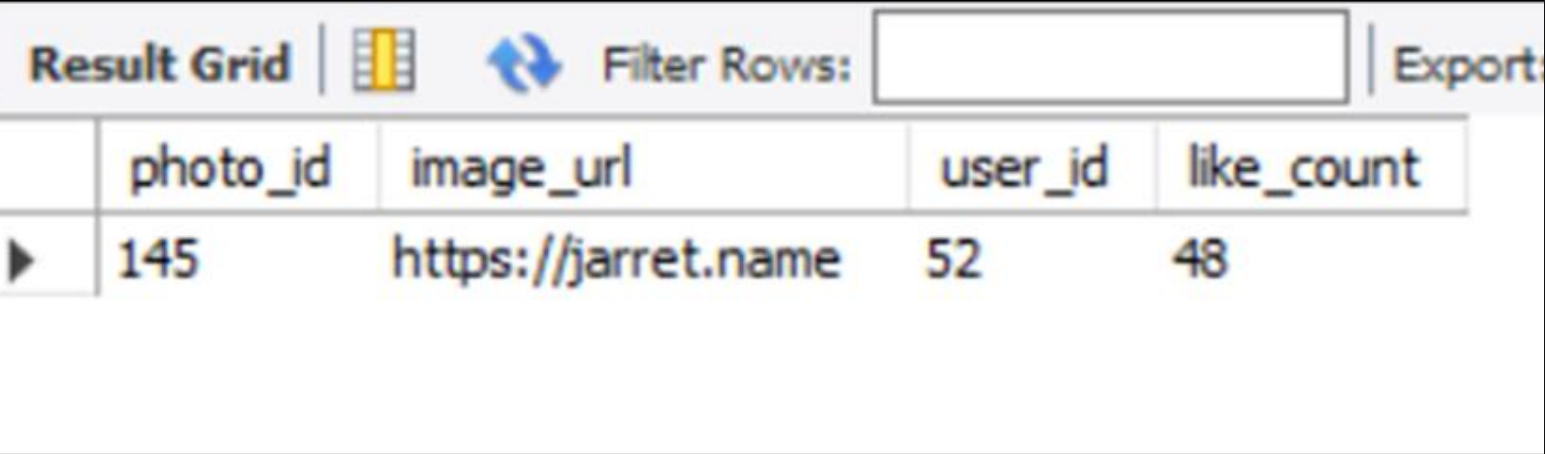This list targets inactive users for promotions.

# Finding – 3

Identify the user with the most likes on a single photo and share their details for the Contest Winner Declaration.



The user details and photo URL will help declare the contest winner.

## Finding - 4

Identify and recommend the top five most used hashtags for the partner brand's outreach.

| tag_name | usage_count |
|----------|-------------|
| smile | 59 |
| beach | 42 |
| party | 39 |
| fun | 38 |
| concert | 24 |

These popular hashtags will be useful for partner brands.

## Finding - 5

Identify the day with the highest user registrations to optimize the ad campaign launch.

| registration_day | user_count |
|------------------|------------|
| Thursday | 16 |

This identifies the day with the highest user registrations.

# Analysis

- Loyal User Reward – Identified the five longest-active users on Instagram to recognize and reward their loyalty, enhancing user retention.
- Inactive User Engagement – Extracted a list of users who have never posted a photo, enabling targeted promotional emails to encourage engagement.
- Contest Winner Declaration – Determined the user with the highest likes on a single photo, ensuring a fair and transparent contest winner selection.
- Hashtag Research – Analyzed and recommended the top five most frequently used hashtags, helping a partner brand maximize reach and engagement.
- Ad Campaign Launch – Evaluated user registration trends to pinpoint the best day for launching ad campaigns, optimizing marketing effectiveness.

## Conclusion

The analysis provides key insights into user engagement, marketing strategies, and investor concerns. By identifying loyal users, inactive accounts, contest winners, and hashtag trends, Instagram can enhance user retention and engagement. Additionally, investor metrics help assess platform authenticity and overall activity levels, guiding future business decisions.

10

# HIRING PROCESS ANALYTICS

## Description

This project analyzes hiring process data to identify patterns, anomalies, and insights for better recruitment decisions. It focuses on handling missing data, aggregating categories, and detecting outliers. Key statistics are summarized to optimize hiring strategies. Data-driven approaches help streamline recruitment and improve decision-making. The goal is to enhance talent acquisition for organizations like Google.

## Design

- Handle Missing Data – Identify and decide how to manage missing values.
- Combine Categories – Merge similar columns for simplified analysis.
- Detect Outliers – Identify data points that may skew results.
- Handle Outliers – Remove, replace, or retain outliers based on impact.
- Summarize Data – Compute key statistics and create visualizations.

# Finding - 1

**Determine the gender distribution of hires by analyzing the number of males and females hired by the company.**



Males led hiring with 2,552 hires, followed by 1,850 females and 277 undisclosed.

# Finding - 2

**Calculate the company's average salary using the Excel formula**



The average salary for named hires was 49,777.70.

# Finding – 3

**Create salary class intervals in Excel using bins**



- Most salaries fall between 40,000 and 50,000.
- Only a small fraction exceed 80,000.

# Finding - 4

**Create a pie chart or bar graph in Excel to visualize department-wise employee distribution.**



- Production led hiring with 39.24%.
- Human Resources followed at 28.34%.
- Sales had the fewest hires at 1.49%.

# Finding – 5

**Create a bar chart or pie chart in Excel to visualize the distribution of employees across different position tiers.**



Distribution of Positions Across Company Tiers

- "C9" was the most assigned position with 1,784 hires.
- "C10" and similar positions had inconsistencies in phrasing.

## Analysis

- Gender Distribution in Hiring: Analyzed the gender ratio among hires to assess diversity and inclusion within the company. This insight can guide future recruitment strategies to ensure balanced representation.
- Salary Trends & Averages: Calculated the average salary offered by the company, providing an understanding of compensation competitiveness and helping HR refine salary benchmarks.
- Salary Distribution Insights: Categorize salaries into class intervals, identifying salary concentration patterns and potential gaps that may need attention for equitable compensation.
- Departmental Workforce Analysis: Visualized department-wise employee distribution, highlighting workforce allocation and identifying departments with high or low hiring rates.

## Conclusion

The hiring process analytics project highlighted how data analysis reveals actionable insights through relationships. Cleaning and visualizing data made patterns and anomalies clearer, guiding recruitment strategies on gender diversity, salary competition, and department-based hiring targets. This project demonstrated how data-driven decisions enhance organizational functions.

# IMDB MOVIE ANALYSIS

## Description

The significance of this study will be exemplified by identifying and investigating IMDB movies' dataset, in order to determine factors affecting the success of a particular movie according to high IMDB ratings. For producers, directors, and business investors, this study has significance in understanding the main elements contributing to success in a movie because they provide for well-grounded decisions when having to start their new projects.

## Design

- Removed missing values and duplicates and ensured standardized formats.
- Fixed issues caused by multiple genre splits.
- Converted columns to appropriate data types.
- Examined the impact of genres, duration, language, director, and budget on IMDB ratings.
- Used descriptive statistics to summarize and compare data.

# Finding - 1

**Identify the most common movie genres in the dataset.**



Drama had the highest number of movies and a high average IMDB score.

# Finding - 2

**Analyze Duration Distribution: Use a histogram in Excel to visualize movie duration ranges**



Most movies had a duration of 90–120 minutes.

# Finding – 3

**Identify Common Languages: Use a bar chart to visualize movie counts per language.**



English dominated with over 80% of the movies.

# Finding – 4

**Identify Top Directors: Calculate each director's average IMDB score**



Akira Kurosawa, Christopher Nolan, and Charles Chaplin had the most frequent and consistently high average IMDB scores

18

# Finding – 5

**Analyze the correlation between movie budgets and gross earnings, and identify the movies with the highest profit margin.**



Top 5 Profit Margin Movies

| | Avatar | Jurassic World | Titanic | Star Wars: Episode IV - A New Hope | E.T. the Extra-Terrestrial |
|---|---|---|---|---|---|
| Series1 | 237000000 | 150000000 | 200000000 | 11000000 | 10500000 |
| Series2 | 523505847 | 502177271 | 458672302 | 449935665 | 424449459 |
| | 760505847 | 652177271 | 658672302 | 460935665 | 434949459 |

Avatar and Jurassic World achieved exceptionally high profit margins relative to their budgets.

## Analysis

- Popular Genres & Ratings – Identified the most common movie genres and analyzed their impact on IMDB scores. Some genres consistently had higher ratings than others.
- Movie Duration & Ratings – Found the average movie duration and its effect on ratings. Longer movies may or may not have higher scores based on trend analysis.
- Language & Ratings – Determined the most common languages in movies and how they affect IMDB ratings. English films showed more rating variations.
- Top Directors & Success – Ranked directors based on their average IMDB scores, highlighting those with the best-performing movies.
- Budget & Profitability – Analyzed the relationship between budgets and earnings, identifying movies with the highest profit margins.

## Conclusion

This project provided valuable insights into the factors influencing IMDB movie ratings, financial success, and industry trends. By analyzing genres, durations, languages, directors, and budgets, we identified key patterns and relationships that impact movie performance. The findings highlight how data-driven decision-making can optimize movie production, marketing, and audience engagement strategies.
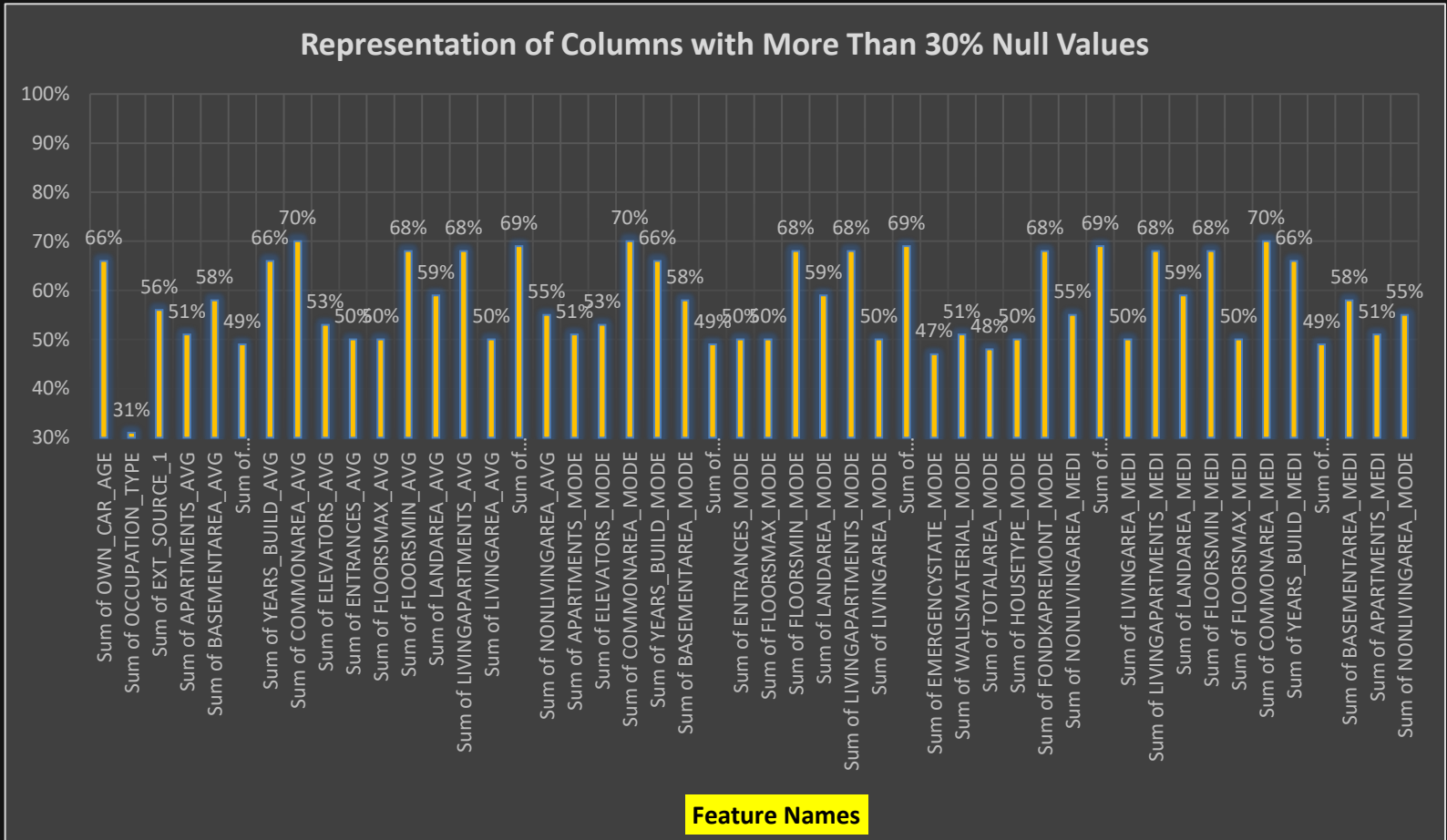
# BANK LOAN CASE STUDY

## Description

The main objective of the project is to analyze the loan application data provided by the customer for identifying patterns that are important when a client applies for any loan about how they are general attributes that raise change in their very default after the first repayment phase. Here the project entails that Explorative Data Analysis must be taken to find out good reasons why a nonpayment or delayed payment has occurred. This will provide useful insights so that they can make data-driven choices for the future.

## Design

- Load and check the dataset for missing values, duplicates, and inconsistencies.
- Handle missing data through imputation or removal while maintaining data integrity.
- Perform single-dimensional analysis to understand feature distributions.
- Conduct bivariate and multivariate analysis to identify correlations and patterns.
- Use bar plots, histograms, scatter plots, and box plots for visualization.
- Compare statistical and visual patterns between defaulters and non-defaulters.

# Finding - 1

**Identify the missing data in the dataset and decide on an appropriate method to deal with it using Excel built-in functions and features.**



COMMONAREA_AVG, NONLIVINGAPARTMENTS_AVG, and FONDKAPREMONT_MODE have over 68% missing values.

# Finding - 2

**Detect and identify outliers in the dataset using Excel statistical functions and features, focusing on numerical variables.**



Income outliers exist in both target groups (0 and 1) but are more pronounced in group 1, suggesting a higher probability of high-income individuals in this group.

# Finding – 3

**Determine if there is data imbalance in the loan application dataset and calculate the ratio of data imbalance using Excel functions.**



- ·Class 0 (No Payment Difficulties): 45,973 instances (91.95% of the dataset).
- ·Class 1 (Payment Difficulties): 4,026 instances (8.05% of the dataset).

# Finding – 4

**Conduct univariate, segmented univariate, and bivariate analysis in Excel to explore variable distributions and relationships with the target variable.**



Default rates are highest among young (0-30) and middle-aged (31-50) borrowers.

# Finding - 5

**Segment the dataset by scenarios (e.g., clients with payment difficulties vs. others) and identify top correlations**



The correlation matrix analyzes relationships among income, credit, demographics, and behavioral flags

## Analysis

- Handling Missing Data – Identified missing values and used Excel functions like AVERAGE or MEDIAN for imputation. Visualized missing data proportions using bar charts.
- Detecting Outliers – Used statistical methods (IQR, QUARTILE) to find extreme values in loan attributes. Box plots and scatter plots helped highlight outliers.
- Assessing Data Imbalance – Checked the distribution of loan approval/rejection cases using COUNTIF. Created pie charts to showcase imbalance in class distribution.
- Univariate & Bivariate Analysis – Analyzed individual variables, segmented data by loan status, and explored relationships using pivot tables, histograms, and scatter plots.
- Finding Key Correlations – Used CORREL function to identify strong indicators of loan default, ranking top correlated factors with heatmaps for different scenarios.

## Conclusion

In conclusion, this project expanded our understanding of banking, particularly in risk assessment and customer analysis. It demonstrated how exploratory data analysis and statistical methods enhance loan approval policies while minimizing risks. Additionally, it highlighted the impact of data imbalances and outliers on decision-making, emphasizing their critical role in banking strategies.

# ANALYZING THE IMPACT OF CAR FEATURES ON PRICE AND PROFITABILITY
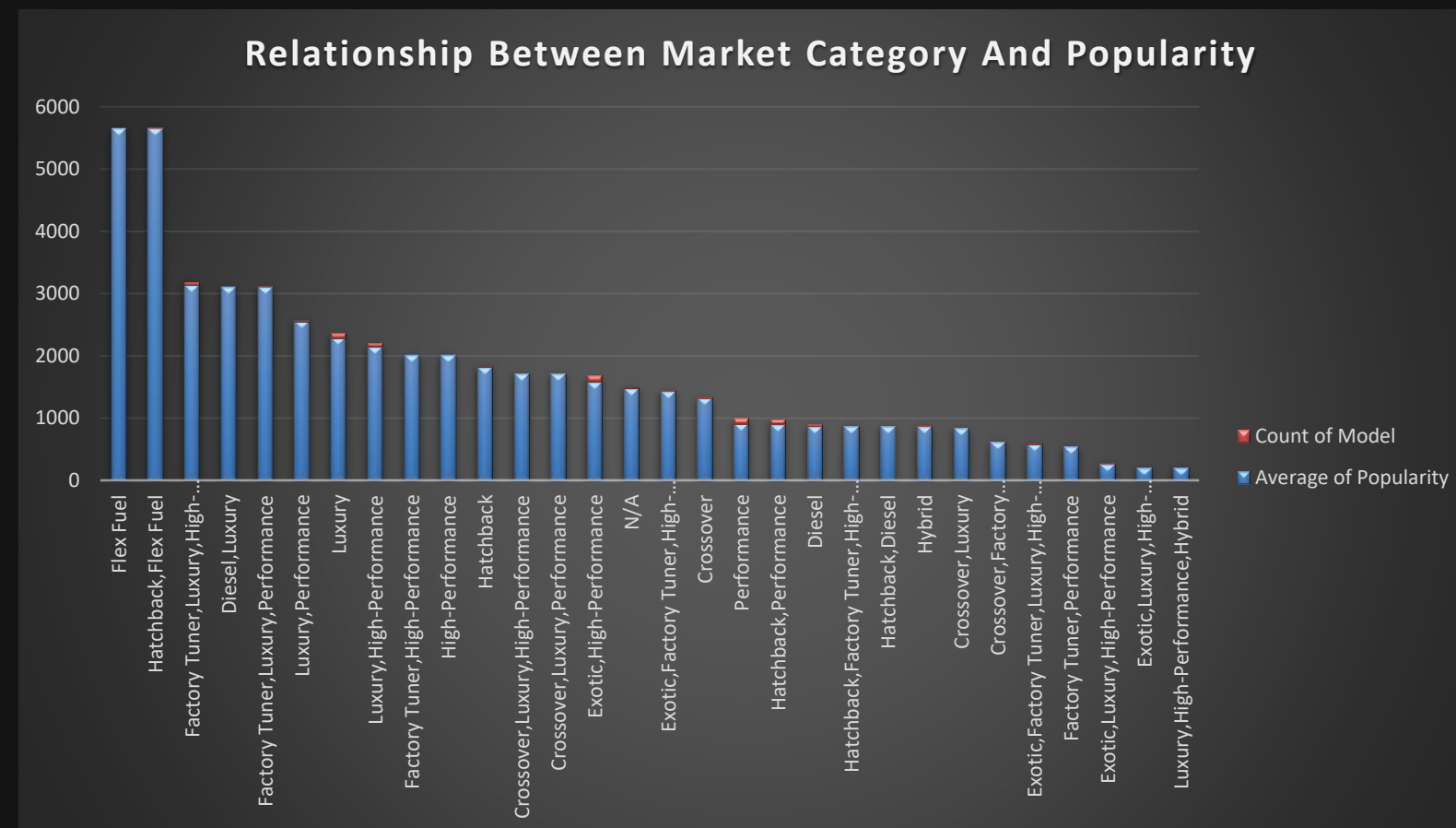
## Description

Change has become very rapid in the automobile industry due to increased consumer interest in fuel efficiency, technological advancement, and environmental sustainability. Automakers would, therefore, need to compare pricing and product development on those fronts that still made business sense to the automakers. This project aims to capture key factors that drive the price of cars and popularity, thereby giving insight in data format on optimizing pricing and product development.

## Design

- Data Cleaning and Preprocessing: Ensured high-quality data for accurate analysis.
- Handling Missing Data: Identified missing values in Engine HP, Engine Cylinders, Market Category, and Number of Doors.
- Used imputation techniques (mean/mode) or removed incomplete data when necessary.
- Data Type Conversion: Converted numerical fields to appropriate data types for consistency.
- Outlier Detection: Identified and managed outliers in MSRP and Engine HP to prevent skewed results.

# Finding - 1

**Create a combo chart that visualizes the relationship between market category and popularity.**



Flex Fuel Tops Popularity – This category has the highest number of models and highest popularity, signifying strong consumer preference

# Finding - 2

**Create a scatter chart with engine power (x-axis) and price (y-axis), adding a trendline to show their relationship.**



There is a general trend showing that as engine HP increases, MSRP (price) also increases, but the relationship is not strictly linear due to significant outliers.

# Finding – 3

**Perform regression analysis to find key variables affecting car price, then visualize coefficient values with a bar chart.**



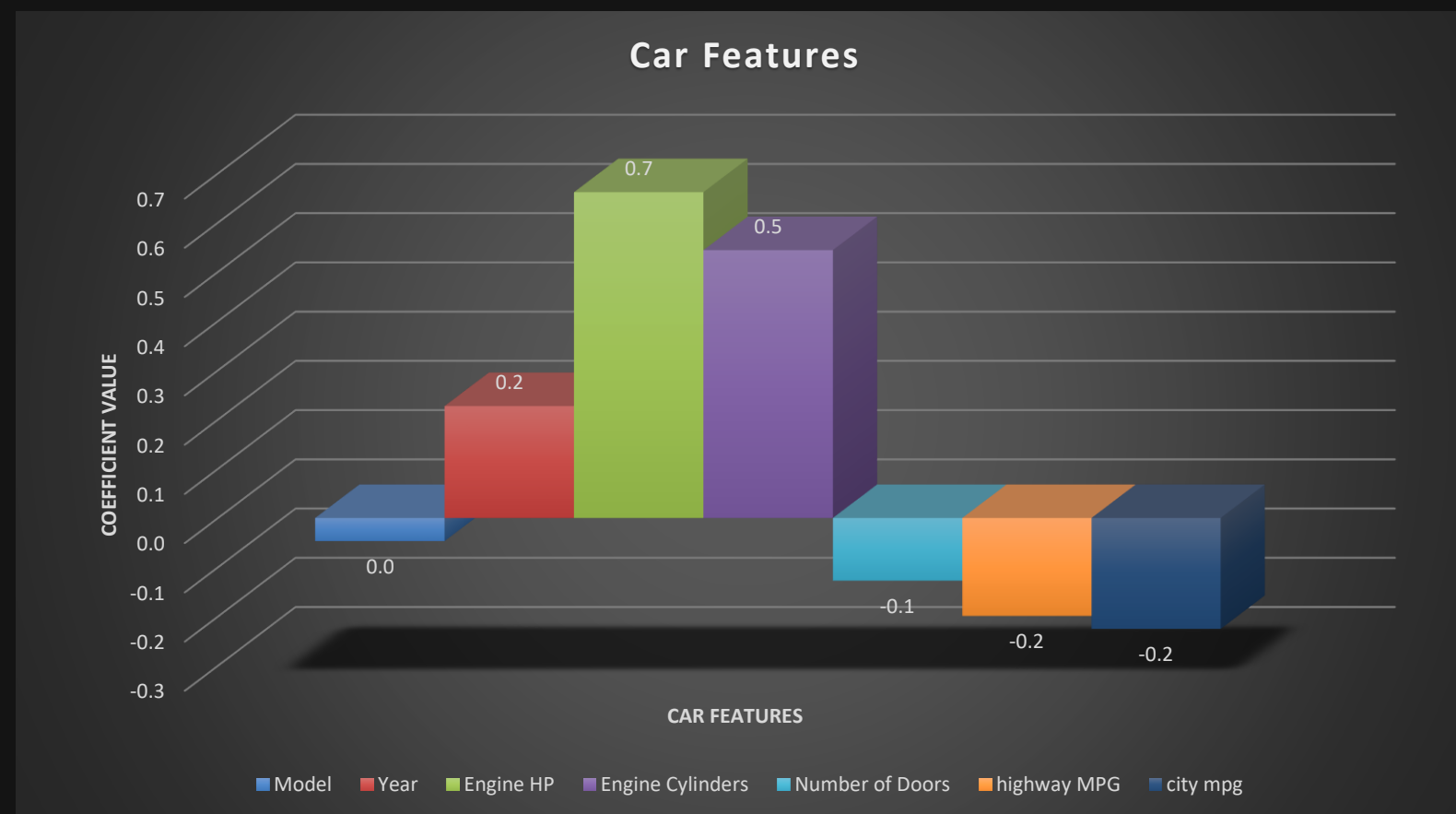Engine HP (0.7) and Cylinders (0.5) have the strongest positive impact, indicating higher power boosts price and performance.

# Finding – 4

**Create a bar chart or a horizontal stacked bar chart that visualizes the relationship between manufacturer and average price.**



Bugatti leads with the highest average price, reflecting its ultra-luxury, high-performance market position.

# Finding – 5

**Create a scatter plot with cylinders (x-axis) and highway MPG (y-axis), adding a trendline to estimate the relationship's slope and significance.**



Relationship between fuel efficiency and the number of cylinders in a Car's engine

More cylinders lead to lower highway MPG, showing that larger engines are less fuel-efficient.

## Analysis

- Car Price Factors: Horsepower and cylinder count greatly impact MSRP; fuel efficiency lowers price.

- Luxury vs. Mainstream: Brands like Bugatti and Ferrari are pricey, while Toyota and Ford target mass markets.

- Fuel Efficiency Over Time: Cars are becoming more fuel-efficient, especially hybrids and smaller engines.

- Data Gaps: Missing engine specs and market categories required imputation.

- Pricing & Feature Challenges: Luxury brands skew pricing, and too many unique categorical features made encoding difficult.

## Conclusion

This analysis helps car manufacturers optimize pricing and product development by identifying profitable features and popular market categories. Using regression analysis and market segmentation, manufacturers can balance consumer demand with profitability. These insights drive strategic pricing and innovation, enhancing competitiveness and long-term success.

# ABC CALL VOLUME TREND ANALYSIS

## Description

Ensuring uninterrupted customer support is key to satisfaction and retention. This project analyzes inbound support call data over 23 days to find insights. The dataset includes agentID, queue time, call duration, and call status. Using data analytics, the goal is to enhance agent efficiency and reduce wait times. Ultimately, this aims to improve the overall customer experience.

## Design

Preprocessing & Data Cleaning:
- Import the dataset and remove missing, duplicate, or inconsistent data.
- Standardize date-time formats and validate data integrity for accurate analysis.

Exploratory Data Analysis (EDA):
- Analyze trends in call volume, call duration, and queue time using statistical methods.
- Use visualizations (bar charts, histograms, time-series plots) to identify patterns in customer calls.
- Detect and assess outliers and their impact on service quality.
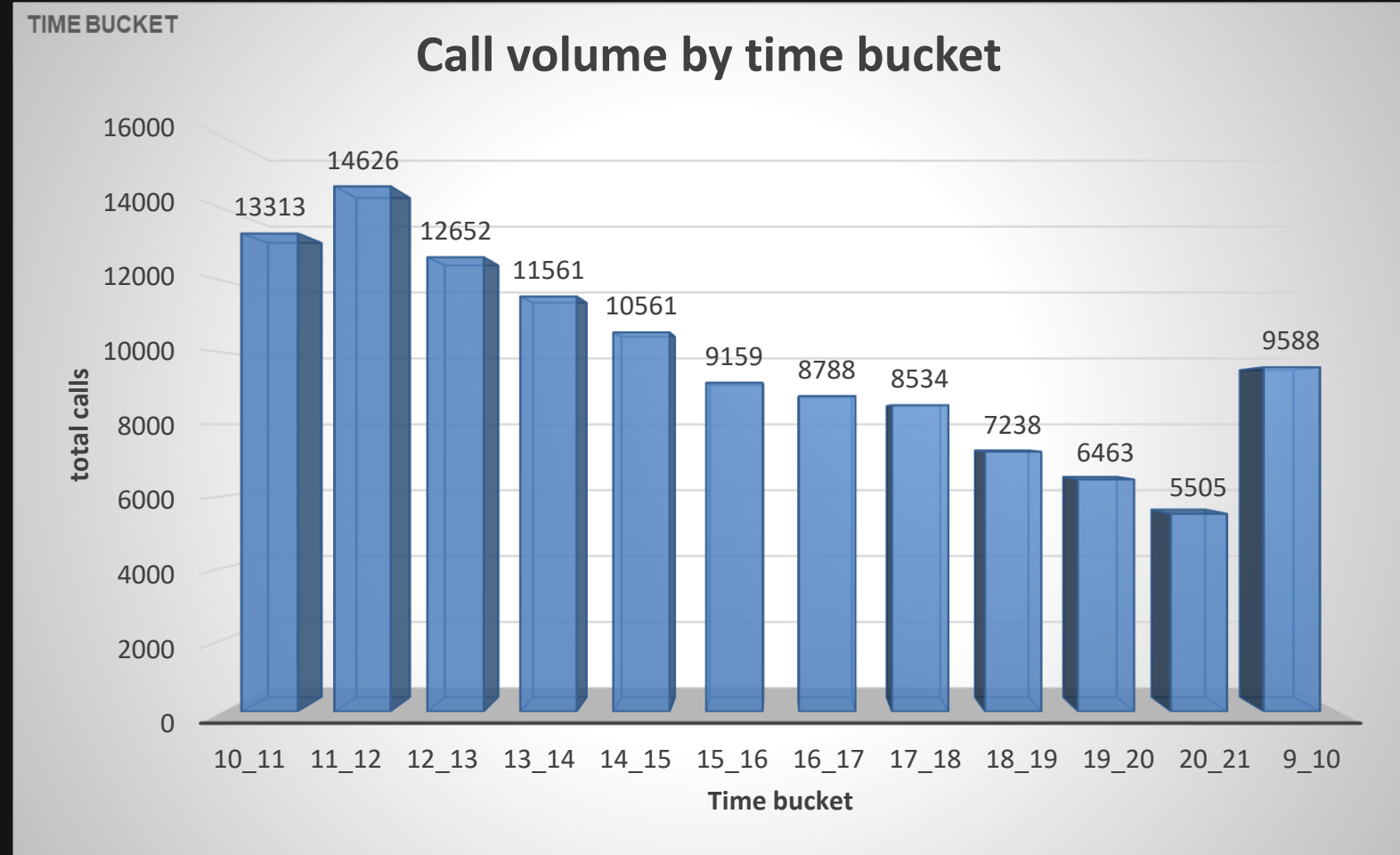
# Finding – 1

**Calculate the average call duration for each time bucket of incoming calls.**



Max Call Duration: 10:00 AM (203.33s) and 7:00 PM (203.41s), likely due to complex queries or engaged agents.
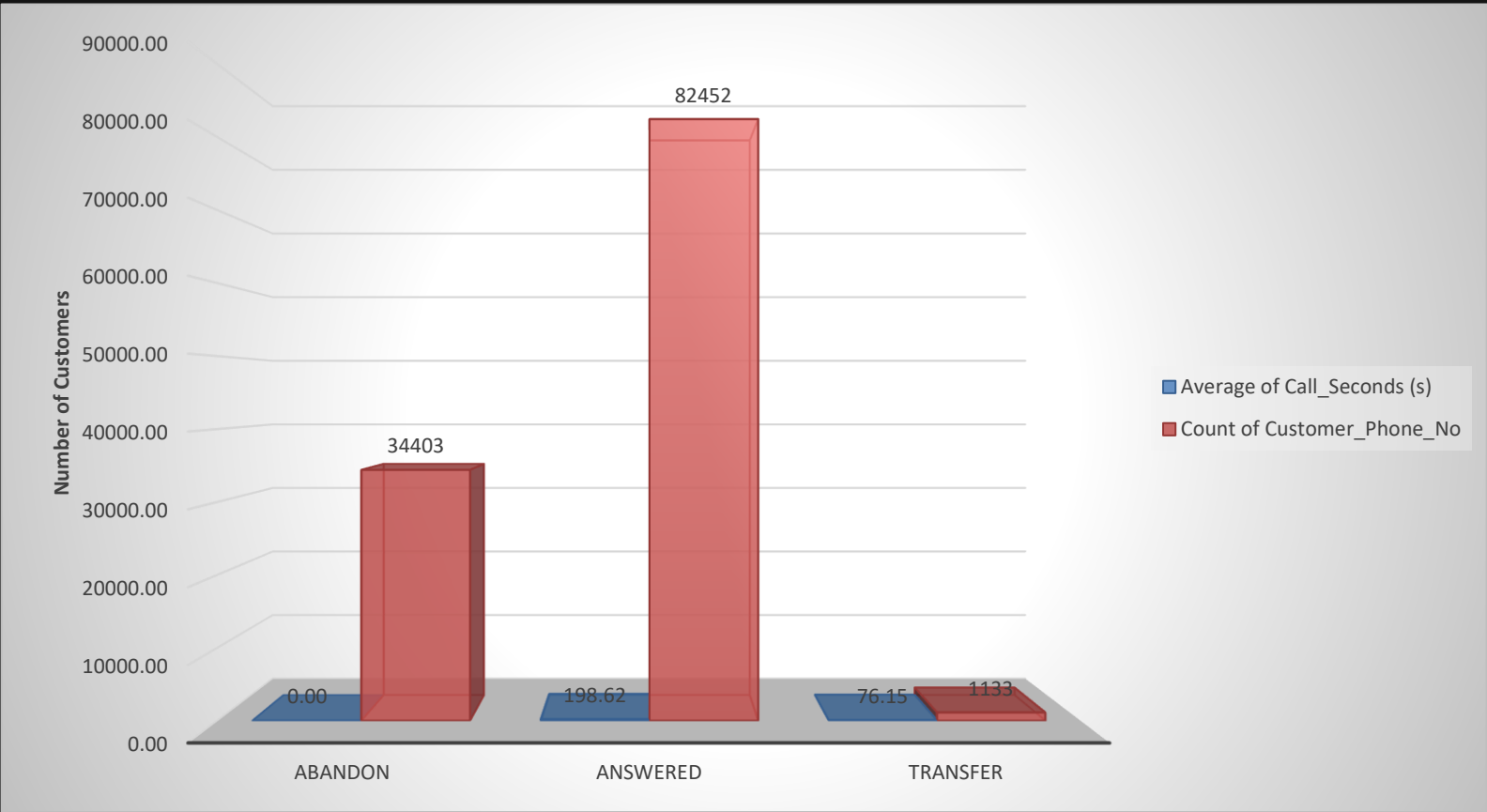
# Finding – 2

**Create a chart showing call volume by time buckets (e.g., 1–2, 2–3).**



Peak call volume: 11–12 PM (14,626 calls) and 10–11 AM (13,313 calls).

# Finding – 3

**Determine the minimum agents needed per time bucket (9 AM – 9 PM) to reduce abandoned calls from 30% to 10%, ensuring 90% answer rate.**



About 30 percent of calls are abandoned, meaning customers are disconnecting before being assisted by an agent.

# Finding – 4

**Design a manpower plan for 24-hour coverage, ensuring a maximum 10% abandon rate, considering 30% of calls occur at night (9 PM – 9 AM).**

| 9PM - 9AM | Call Distribution | Time Distribution | Agent Required |
|---|---|---|---|
| 9_10 | 3 | 10 | 1.50 |
| 10_11 | 3 | 10 | 1.50 |
| 11_12 | 2 | 15 | 1.00 |
| 12_1 | 2 | 15 | 1.00 |
| 1_2 | 1 | 30 | 0.50 |
| 2_3 | 1 | 30 | 0.50 |
| 3_4 | 1 | 30 | 0.50 |
| 4_5 | 1 | 30 | 0.50 |
| 5_6 | 3 | 10 | 1.50 |
| 6_7 | 4 | 7.5 | 2.00 |
| 7_8 | 4 | 7.5 | 2.00 |
| 8_9 | 5 | 6 | 2.50 |
| Total | 30 | | 15.00 |

| Average daily call | 5130 |
|---|---|
| For Night | 1539 |
| Additional Hour Required | 76 |
| Additional Agent Required | 15 |

Unanswered Night Calls- No Agents from 9pm to 9 am, which leads to poor customer experience

## Analysis

- Call Volume Trends: Peak hours are 10 AM–1 PM, with the highest at 11 AM–12 PM (14,626 calls). Calls drop sharply after 3 PM, with minimal activity after 8 PM.
- Call Duration Analysis: Longest calls occur at 10–11 AM and 7–8 PM (~203 sec), while shortest calls are between 12–3 PM, indicating quick resolutions or shorter wait times.
- Call Abandonment & Queue Times: 30% of calls are abandoned due to long wait times; reducing this by 10% requires better staffing and handling strategies.
- Agent Staffing Needs: Peak hours (10 AM–2 PM, 5–6 PM) need 6–7 agents, while low-volume hours (7 PM–9 AM) need 3–4. A total of 56 agents ensures a 90% answer rate.
- Night Shift Planning: No agents between 9 PM–9 AM leads to poor service. To handle 1,539 nightly calls, 15 more agents are needed, with demand peaking from 6–9 AM.

## Conclusion

This project analyzed call volumes, agent performance, and workforce optimization to enhance customer experience and operational efficiency. Identifying peak hours reduced abandonment rates and improved resource allocation. AI-driven solutions like IVR and chatbots boosted self-service options, achieving 90% call-handling efficiency and seamless 24/7 customer support.

# APPENDIX

1. **Instagram User Analytics.** ----------- [Click Here](#)

2. **Hiring Process Analytics.** ------------- [Click Here](#)

3. **IMDB Movie Analysis.** ---------------- [Click Here](#)

4. **Bank Loan Case Study.** --------------- [Click Here](#)

5. **Impact of Car Features.** -------------- [Click Here](#)

6. **ABC Call Volume Trend Analysis.** ------- [Click Here](#)

# CONTACT INFO

## Email Address

mohan.s.agnivamsha@gmail.com