

Data Analysis Portfolio

Jeevan Kishore



Professional Background

I am Jeevan Kishore, I have Completed B.E. Electronics and Communication Engineering in 2022 from K.Ramakrishnan College of Technology, Trichy. I have secured 8.01 CGPA and have several skills including Data Analysis, Python, MySQL, Excel, Tableau, R Programming.

I'm a quick learner and passionate about staying up-to-date with the latest trends and technologies. With my diverse skill set, I am capable of working on a wide range of projects. I have excellent communication and collaboration skills and enjoy working in a team environment to drive project success.

Table of Contents

| S. No. | Description | Page No. |
|--------|---------------------------------------------------|----------|
| 1 | Professional Background | 2 |
| 2 | Table of Contents | 3 |
| 3 | Data Analysis Process | 4-6 |
| 4 | Instagram User Analytics | 7-16 |
| 5 | Operation and Metric Analytics | 17-23 |
| 6 | Hiring Process Analysis | 24-30 |
| 7 | IMDb Movie Analysis | 31-39 |
| 8 | Bank Loan Case Study | 40-61 |
| 9 | Impact of Car Features on Price and Profitability | 62-77 |
| 10 | ABC Call Volume Trend | 78-86 |
| 11 | Drive Links | 87 |

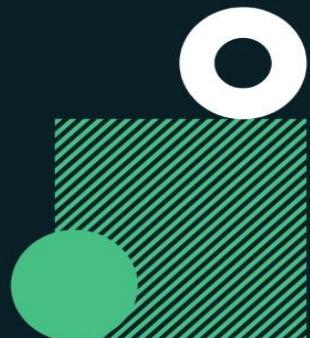
Data Analytics Process

Project-1

trainity

Data Analytics Process: Real World Application

Shopping & Use of 6 Step Data Analytics Process



Project Description:

We use Data Analytics in everyday life without even knowing it. Your task is to give the example(s) of such a real-life situation where we use Data Analytics and link it with the data analytics process.

Data Analytics Process

1. Scenario: Joining data analytics course in the best learning platform
 - PLAN : Planned to shift my career to data analytics domain.
 - PREPARE : Then decided to join a course based on my budget and time that I can devote for learning.
 - PROCESS : Then I searched and shortlisted few renowned learning platforms like Trainity, Udemy, Upgrad, etc..
 - ANALYZE : By analyzing these platforms through various criteria like How it works?, What features they offers?, Placement Track, Fees Structure, etc.. and finalized the best platform.
 - SHARE : I registered for a webinar with trainity where I shared my willingness to enroll in course.
 - ACT : I paid the fees and enrolled in course and started learning.

2. Scenario: Cooking a dish

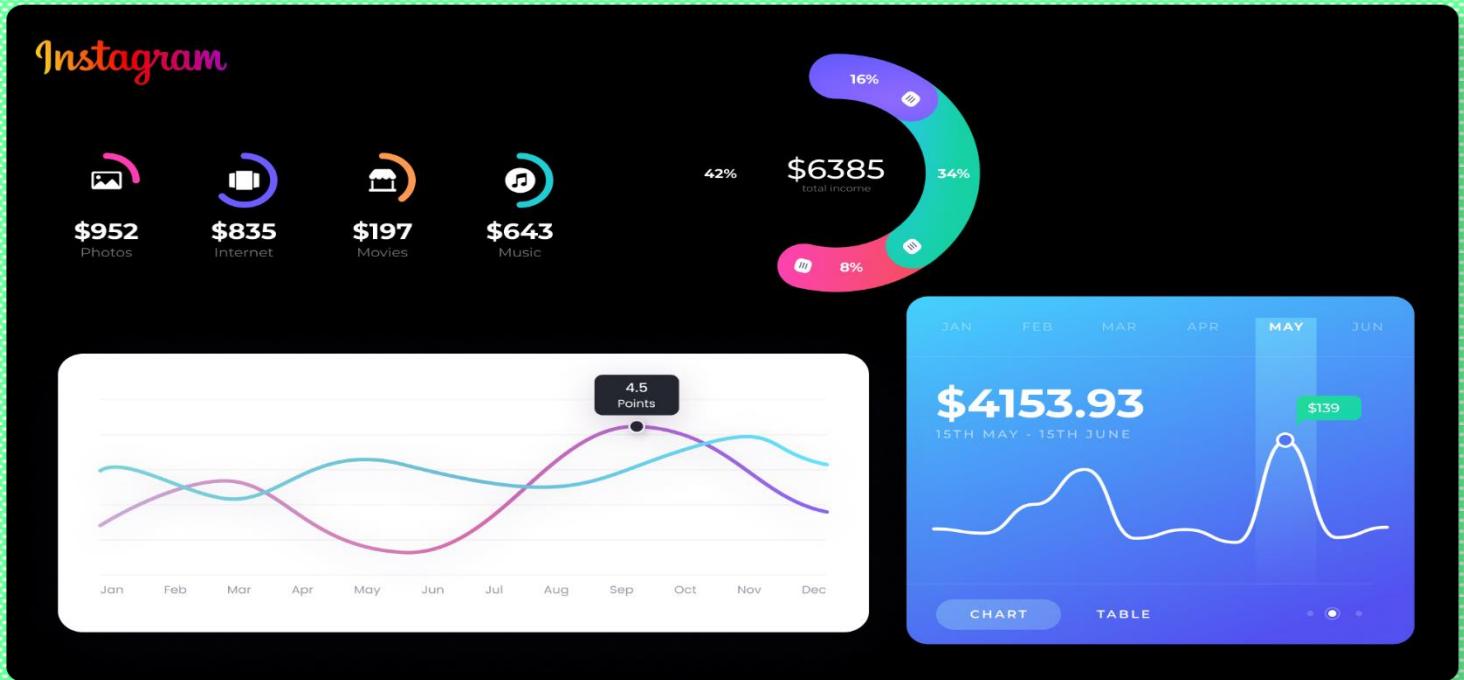
- PLAN : Planning to make non veg for lunch.
- PREPARE : Decided to cook chicken biryani based on the preferences and time constrain to cook.
- PROCESS : Went to shop to purchase chicken, rice, curry leaves, garlic, masalas, etc..
- ANALYZE : Ratio of rice, water quantity to added, ratio of meat, etc.. and the appropriate cooking time for chicken to be cooked.
- SHARE : Conveyed to my friends that I going to cook chicken biryani.
- ACT : The ingredients were organized and started cooking.

Conclusion

Hence, we have seen how we can use the 6 steps of Data Analytics Process while making any decision in real life scenarios.

Instagram User Analytics

Project-2



Project Description :

The aim of this project is to perform analysis on Instagram user data and gain insights into the behavior of Instagram users. SQL was used as a primary tool for data analysis and retrieval. The findings from the analysis provide valuable insights into user activity, preferences and demographics on the popular social media platform. This report presents the methodology and results of the analysis and discusses the implications of the findings for businesses and marketers.

Approach:

Created a database using MySQL that can store the data and support the queries that needs to be processed for analyzing and loaded the data into the database and used MySQL queries to perform the analysis and made insights about the queries that were asked by the management team.

Tech-Stack Used:

MySQL - To run the Query and extract data from database.

MS Excel – To Store the Extracted tables from the database.

MS PPT – To prepare the report based the insights from the database.

Findings

A) Marketing:

1. Rewarding Most Loyal Users:

The screenshot shows a MySQL Workbench interface. At the top, there's a toolbar with various icons. Below it is a text area containing a SQL query:

```
1
2      ##### Rewarding Most Loyal Users
3
4 •  use ig_clone;
5 •  select * from users
6    order by date (created_at)
7    limit 5;
8
```

Below the query is a "Result Grid" table with the following data:

| | id | username | created_at |
|---|------|------------------|---------------------|
| ▶ | 80 | Darby_Herzog | 2016-05-06 00:14:21 |
| | 67 | Emilio_Bernier52 | 2016-05-06 13:04:30 |
| | 63 | Elenor88 | 2016-05-08 01:30:41 |
| | 95 | Nicole71 | 2016-05-09 17:30:22 |
| * | 71 | Nia_Haag | 2016-05-14 15:38:50 |
| | HULL | NULL | NULL |

They are the 5 oldest users of the Instagram from the database provided.

2. Remind Inactive Users to Start Posting:

```
8
9      ### Remind Inactive Users to Start Posting
10 •  select * from users
11    left join photos
12    on users.id = photos.user_id
13    where user_id is null;
14
```

| Result Grid | | | | | | | |
|-------------|----|--------------------|---------------------|------|-----------|---------|-------------|
| | id | username | created_at | id | image_url | user_id | created_dat |
| ▶ | 5 | Aniya_Hackett | 2016-12-07 01:04:39 | NULL | NULL | NULL | NULL |
| | 7 | Kassandra_Homenick | 2016-12-12 06:50:08 | NULL | NULL | NULL | NULL |
| | 14 | Jadyn81 | 2017-02-06 23:29:16 | NULL | NULL | NULL | NULL |
| | 21 | Rocio33 | 2017-01-23 11:51:15 | NULL | NULL | NULL | NULL |
| | 24 | Maxwell.Halvorson | 2017-04-18 02:32:44 | NULL | NULL | NULL | NULL |
| | 25 | Tierra.Trantow | 2016-10-03 12:49:21 | NULL | NULL | NULL | NULL |
| | 34 | Pearl7 | 2016-07-08 21:42:01 | NULL | NULL | NULL | NULL |
| | 36 | Ollie_Ledner37 | 2016-08-04 15:42:20 | NULL | NULL | NULL | NULL |
| | 41 | Mckenna17 | 2016-07-17 17:25:45 | NULL | NULL | NULL | NULL |
| | 45 | David.Osinski47 | 2017-02-05 21:23:37 | NULL | NULL | NULL | NULL |
| | 49 | Morgan.Kassulke | 2016-10-30 12:42:31 | NULL | NULL | NULL | NULL |
| | 53 | Linnea59 | 2017-02-07 07:49:34 | NULL | NULL | NULL | NULL |
| | 54 | Duane60 | 2016-12-21 04:43:38 | NULL | NULL | NULL | NULL |
| | 57 | Julien_Schmidt | 2017-02-02 23:12:48 | NULL | NULL | NULL | NULL |
| | 66 | Mike_Auer39 | 2016-07-01 17:36:15 | NULL | NULL | NULL | NULL |

Result 5 ×

```
8
9      ### Remind Inactive Users to Start Posting
10 •  select * from users
11    left join photos
12    on users.id = photos.user_id
13    where user_id is null;
```

| Result Grid | | | | | | | |
|-------------|----|---------------------|---------------------|------|-----------|---------|-------------|
| | id | username | created_at | id | image_url | user_id | created_dat |
| | 54 | Duane60 | 2016-12-21 04:43:38 | NULL | NULL | NULL | NULL |
| | 57 | Julien_Schmidt | 2017-02-02 23:12:48 | NULL | NULL | NULL | NULL |
| | 66 | Mike_Auer39 | 2016-07-01 17:36:15 | NULL | NULL | NULL | NULL |
| | 68 | Franco_Keebler64 | 2016-11-13 20:09:27 | NULL | NULL | NULL | NULL |
| | 71 | Nia_Haag | 2016-05-14 15:38:50 | NULL | NULL | NULL | NULL |
| | 74 | Hulda.Macejkovic | 2017-01-25 17:17:28 | NULL | NULL | NULL | NULL |
| | 75 | Leslie67 | 2016-09-21 05:14:01 | NULL | NULL | NULL | NULL |
| | 76 | Janelle.Nikolaus81 | 2016-07-21 09:26:09 | NULL | NULL | NULL | NULL |
| | 80 | Darby_Herzog | 2016-05-06 00:14:21 | NULL | NULL | NULL | NULL |
| | 81 | Esther.Zulauf61 | 2017-01-14 17:02:34 | NULL | NULL | NULL | NULL |
| | 83 | Bartholome.Bernhard | 2016-11-06 02:31:23 | NULL | NULL | NULL | NULL |
| | 89 | Jessyca_West | 2016-09-14 23:47:05 | NULL | NULL | NULL | NULL |
| | 90 | Esmeralda.Mraz57 | 2017-03-03 11:52:27 | NULL | NULL | NULL | NULL |
| | 91 | Bethany20 | 2016-06-03 23:31:53 | NULL | NULL | NULL | NULL |

3. Declaring Contest Winner:

The screenshot shows a MySQL Workbench interface. At the top, there's a toolbar with various icons. Below it is a SQL editor window containing the following code:

```
15  ### Declaring Contest Winner
16 • select users.username, users.id, count(likes.photo_id) as total_likes
17  from users
18  left join photos
19  on users.id = photos.user_id
20  left join likes
21  on photos.id = likes.photo_id
22  group by users.username, likes.photo_id, users.id
23  order by count(likes.photo_id) desc
24  limit 1;
25
```

Below the SQL editor is a results grid labeled "Result Grid". It has columns for "username", "id", and "total_likes". A single row is displayed, showing "Zack_Kemmer93" in the "username" column, "52" in the "id" column, and "48" in the "total_likes" column.

Zack_Kemmer93 is the winner with most number of likes for a single photo.

4. Hashtag Researching :

```
25
26      ### Top 5 Hashtag Researching
27 •   select tag_name as "TOP 5 Tags", count(tags.id) as "No. of Times Used"
28   from photos
29   left join photo_tags
30   on photos.id = photo_tags.photo_id
31   left join tags
32   on photo_tags.tag_id = tags.id
33   Group by tags.id
34   order by count(id) desc
35   limit 5;
~~
```

| Result Grid | | |
|-------------|------------|-------------------|
| | TOP 5 Tags | No. of Times Used |
| ▶ | smile | 59 |
| | beach | 42 |
| | party | 39 |
| | fun | 38 |
| | concert | 24 |

| tag_name | count(t.id) |
|----------|-------------|
| smile | 59 |
| beach | 42 |
| party | 39 |
| fun | 38 |
| concert | 24 |

The list of all the hashtags were the most used in Instagram.

5. Launch AD Campaign:

```
36
37      ### Best Day for Launch of AD Campaign
38 •  select dayname(created_at) as Days_of_the_week, count(*) as Total
39      from users
40      group by Days_of_the_week
41      order by Total desc
42
43
44
```

| Result Grid | Filter Rows: | Export: | Wrap Cell Content:

| | Days_of_the_week | Total |
|---|------------------|-------|
| ▶ | Thursday | 16 |
| | Sunday | 16 |
| | Friday | 15 |
| | Tuesday | 14 |
| | Monday | 14 |
| | Wednesday | 13 |
| | Saturday | 12 |

Result 34 ×

| days_of_week | count(days_of_week) |
|--------------|---------------------|
| Thursday | 16 |
| Sunday | 16 |
| Tuesday | 14 |
| Saturday | 12 |
| Wednesday | 13 |
| Monday | 14 |
| Friday | 15 |

These days of the week were the most users registered on.

B) Investor Metrics

6. User Engagement:

The screenshot shows the MySQL Workbench interface. The SQL editor contains the following code:

```
43     ###Investor Metrics
44     ###User Engagement
45
46     #total no of users of instagram
47 •   select count(id) as total_no_of_users from users;
```

The result grid shows one row with the value 100 under the column 'total_no_of_users'.

Total Number of Users in Instagram

The screenshot shows the MySQL Workbench interface. The SQL editor contains the following code:

```
43     ###Investor Metrics
44     ###User Engagement
45
46     #total no of users of instagram
47 •   select count(id) as total_no_of_users from users;
48
49     #total no of photos on instagram
50 •   select count(id) as total_no_of_photos from photos;
```

The result grid shows one row with the value 257 under the column 'total_no_of_photos'.

Total Number of Photos in Instagram

```

51
52      #average post per user
53 •  with cte as
54   ( select users.id, count(photos.id) as post_per_user
55     from users
56     left join photos
57       on users.id = photos.user_id
58     group by users.id
59     order by post_per_user desc )
60   select avg(post_per_user) as "Average Post per User"
61   from cte;
~~

```

| Result Grid | | Filter Rows: | Export: | Wrap Cell Content: |
|-----------------------|--|--------------|---------|-----------------------|
| Average Post per User | | | | |
| ▶ 2.5700 | | | | Average Post per User |

7. Bots & Fake Accounts:

```

63      #####Bots & Fake Accounts
64 •  with cte as
65   (select likes.user_id, count(likes.user_id) as no_of_likes
66     from photos
67     left join likes
68       on photos.id = likes.photo_id
69     left join users
70       on photos.user_id = users.id
71     group by likes.user_id)
72   select user_id, users.username,no_of_likes
73   from cte
74   left join users
75     on cte.user_id = users.id
76   where no_of_likes = 257;
77

```

| Result Grid | | Filter Rows: | Export: | Wrap Cell Content: |
|-------------|--------------------|--------------|---------|--------------------|
| user_id | username | no_of_likes | | |
| ▶ 5 | Aniya_Hackett | 257 | | |
| 14 | Jadyn81 | 257 | | |
| 21 | Rocio33 | 257 | | |
| 24 | Maxwell.Halvorson | 257 | | |
| 36 | Ollie_Ledner37 | 257 | | |
| 41 | Mckenna17 | 257 | | |
| 54 | Duane60 | 257 | | |
| 57 | Julien_Schmidt | 257 | | |
| 66 | Mike.Auer39 | 257 | | |
| 71 | Nia_Haag | 257 | | |
| 75 | Leslie67 | 257 | | |
| 76 | Janelle.Nikolaus81 | 257 | | |
| 91 | Bethany20 | 257 | | |

Result 55 ×

These are the accounts that have liked all the photos that are available, Hence we can conclude that these are bots.

Conclusion

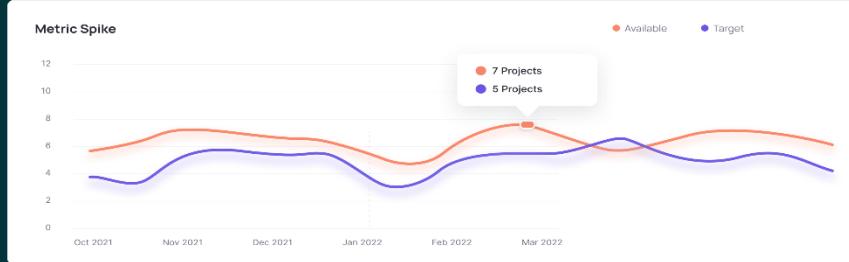
Analyzed many useful insights that could help the business to launch a new marketing campaign, decide on features to build for an app, track the success of the app by measuring user engagement and improve the experience altogether would help the business grow and this would help in making data driven decisions for the business.

Operational and Metrics Analytics

Project-3

trainity

Operation Analytics & Investigating metric spike case study



Project Description:

This project involved using SQL to analyse operational data and investigate metric spikes. The goal was to generate insights and answer complex questions related to the organization's operations. I used a range of SQL techniques to answer complex questions, such as subqueries, joins, and aggregations. For example, I used subqueries to compare the performance of different teams or regions, and joins to combine data from multiple

sources. I also calculated metrics such as rolling averages and percentiles to gain a deeper understanding of the data.

Approach:

Use SQL to create a database that can store the data and support the queries that we need to perform the analysis and load the data into the database. Use SQL queries to perform the analysis and retrieve the insights you want to extract.

Tech-Stack Used:

MySQL Workbench 8.0 CE

Case Study 1 (Job Data)

1. Number of jobs reviewed:

The screenshot shows the MySQL Workbench interface. At the top, there's a toolbar with various icons. Below the toolbar, the SQL editor contains the following code:

```
1 •  use cs1;
2
3      ###Number of jobs reviewed
4 •  select ds, count(job_id) as 'Jobs Reviewed per day', sum(time_spent)
5      from job_data
6      where ds >='2020-11-01'  and ds <='2020-11-30'
7      group by ds ;
```

Below the SQL editor is the Result Grid pane, which displays the query results as a table:

| | ds | Jobs Reviewed per day | Hours spent per Day |
|---|------------|-----------------------|---------------------|
| ▶ | 2020-11-30 | 2 | 0.0111 |
| | 2020-11-29 | 1 | 0.0056 |
| | 2020-11-28 | 2 | 0.0092 |
| | 2020-11-27 | 1 | 0.0289 |
| | 2020-11-26 | 1 | 0.0156 |
| | 2020-11-25 | 1 | 0.0125 |

The Result Grid pane has several buttons at the top: 'Result Grid', 'Filter Rows:', 'Export:', and 'Wrap Cell Content:'. The status bar at the bottom left shows 'Result 3'.

2.Throughput:

```
9      #####Throughput
10 •   SELECT ds,
11      sum(jobs) over (order by ds rows between 6 preceding and current row) /
12      sum(total_time) as throughput_7day_rolling_avg
13      from
14      (SELECT ds, COUNT(job_id) as jobs, SUM(time_spent) as total_time
15       FROM job_data where ds >= '2020-11-01' and ds <='2020-11-30'
16       GROUP BY ds
17      ) d
18      group by ds;
```

| Result Grid | | Filter Rows: | Export: | Wrap Cell Content: |
|-------------|------------|---------------------------|---------|--------------------|
| | ds | throughput_7d_rolling_avg | | |
| ▶ | 2020-11-25 | 0.0222 | | |
| | 2020-11-26 | 0.0357 | | |
| | 2020-11-27 | 0.0288 | | |
| | 2020-11-28 | 0.1515 | 0.1515 | |
| | 2020-11-29 | 0.3000 | | |
| | 2020-11-30 | 0.2000 | | |

3.Percentage share of each language:

```
19
20      #####Percentage share of each language
21 •   select language, (count(*)/(select count(*) from job_data))*100 as percentage
22     from job_data
23     group by language;
```

| Result Grid | | Filter Rows: | Export: | Wrap Cell Content: |
|-------------|----------|--------------|---------|--------------------|
| | language | percentage | | |
| ▶ | English | 12.5000 | | |
| | Arabic | 12.5000 | | |
| | Persian | 37.5000 | | |
| | Hindi | 12.5000 | | |
| | French | 12.5000 | | |
| | Italian | 12.5000 | | |

4.Duplicate rows:

```
24
25      ###Duplicate rows
26 •  select job_id, count(job_id) as duplicate_rows
27      from job_data
28      group by job_id
29      having count(*) > 1;
30
31
32
33
34
```

| Result Grid | | |
|-------------|--------|----------------|
| | job_id | duplicate_rows |
| ▶ | 23 | 3 |

Case Study 2 (Investigating metric spike)

1.User Engagement:

```
3 •  SELECT
4      extract(year from occurred_at) as year,
5      extract(week from occurred_at) as weeknum,
6      count(distinct user_id) as user_enagement
7  FROM
8      events
9  GROUP BY
10     year, weeknum
11  order by
12     year, weeknum
```

| Result Grid | | | |
|-------------|------|---------|----------------|
| | year | weeknum | user_enagement |
| ▶ | 2014 | 17 | 663 |
| | 2014 | 18 | 1068 |
| | 2014 | 19 | 1113 |
| | 2014 | 20 | 1154 |
| | 2014 | 21 | 1121 |

2. User Growth:

```
3 • select
4   year_num,
5   week_num,
6   num_active_users,
7   SUM(num_active_users)OVER(ORDER BY year_num, week_num ROWS BETWEEN
8   UNBOUNDED PRECEDING AND CURRENT ROW) AS cum_active_users
9   from
10  (
11    select
12      extract(year from activated_at) as year_num,
13      extract(week from activated_at) as week_num,
14      count(distinct user_id) as num_active_users
15    from
16    users
17    WHERE
18      state = 'active'
19    group by year_num, week_num
20    order by year_num, week_num
21  ) a;
~~
```

| Result Grid | | | |
|-------------|----------|------------------|------------------|
| year_num | week_num | num_active_users | cum_active_users |
| 2013 | 0 | 23 | 23 |
| 2013 | 1 | 30 | 53 |
| 2013 | 2 | 48 | 101 |
| 2013 | 3 | 36 | 137 |
| 2013 | 4 | 30 | 167 |
| 2013 | 5 | 48 | 215 |
| 2013 | 6 | 38 | 253 |
| 2013 | 7 | 47 | 295 |

3. Weekly Retention:

```
23 • SELECT
24   distinct user_id,
25   COUNT(user_id),
26   SUM(CASE WHEN retention_week = 1 Then 1 Else 0 END) as per_week_retention
27   FROM
28  (
29    SELECT
30      a.user_id,
31      a.signup_week,
32      b.engagement_week,
33      b.engagement_week - a.signup_week as retention_week
34    FROM
35    (
36      (SELECT distinct user_id, extract(week from occurred_at) as signup_week FROM events_data
37      WHERE event_type = 'signup_flow'
38      and event_name = 'complete_signup'
39      and extract(week from occurred_at) = 18
40      )a
41      LEFT JOIN
42      (SELECT distinct user_id, extract(week from occurred_at) as engagement_week FROM events_data
43      where event_type = 'engagement'
```

| Result Grid | | |
|-------------|----------------|--------------------|
| user_id | COUNT(user_id) | per_week_retention |
| 11926 | 1 | 0 |
| 11928 | 1 | 0 |
| 11929 | 1 | 0 |
| 11931 | 1 | 0 |
| 11933 | 1 | 0 |

```

29   SELECT
30     a.user_id,
31     a.signup_week,
32     b.engagement_week,
33     b.engagement_week - a.signup_week AS retention_week
34   FROM
35   (
36     (SELECT DISTINCT user_id, EXTRACT(week FROM occurred_at) AS signup_week FROM events_data
37      WHERE event_type = 'signup_flow'
38      AND event_name = 'complete_signup'
39      AND EXTRACT(week FROM occurred_at) = 18
40    )a
41   LEFT JOIN
42     (SELECT DISTINCT user_id, EXTRACT(week FROM occurred_at) AS engagement_week FROM events_data
43      WHERE event_type = 'engagement'
44    )b
45   ON a.user_id = b.user_id
46   )
47   GROUP BY user_id
48   ORDER BY user_id;
49

```

| Result Grid | | |
|-------------|---------|----------------|
| | | |
| | user_id | COUNT(user_id) |
| | 11936 | 1 |
| | 11939 | 3 |
| | 11940 | 1 |
| | 11942 | 1 |
| | 11944 | 1 |

4. Weekly Engagement:

```

50
51 •  SELECT
52   EXTRACT(YEAR FROM occurred_at) AS year_num,
53   EXTRACT(WEEK FROM occurred_at) AS week_num,
54   device,
55   COUNT(DISTINCT user_id) AS no_of_users
56   FROM
57   events_data
58   WHERE event_type = 'engagement'
59   GROUP BY 1,2,3
60   ORDER BY 1,2,3;
61
62
63
64

```

| Result Grid | | | |
|-------------|----------|----------|-----------------------|
| | year_num | week_num | device |
| ▶ | 2014 | 17 | acer aspire desktop |
| | 2014 | 17 | acer aspire notebook |
| | 2014 | 17 | amazon fire phone |
| | 2014 | 17 | asus chromebook |
| | 2014 | 17 | dell inspiron desktop |

5.Email Engagement:

```
61
62 •   SELECT
63     engagements,
64     total_users,
65     engagements/total_users*100 AS engagement_rate
66   FROM (
67       SELECT
68           COUNT(DISTINCT
69           CASE
70             WHEN action = 'email_open' THEN user_id
71             WHEN action = 'email_clickthrough' THEN user_id END) AS engagements
72         FROM email_events) AS counts;
73
74
75
```

| Result Grid | | | |
|-------------|-------------|-------------|-----------------|
| | engagements | total_users | engagement_rate |
| ▶ | 5927 | 6179 | 95.9217 |

Conclusion

Extracted many useful analysis that could help in making data driven decisions making for the business.

HIRING PROCESS ANALYTICS

Project-4

Project Description:

The project aims is to analyze the hiring process of a company and provide a report about the analysis made out of the data provided. The analysis has been done using Excel and statistical techniques.

Approach:

The dataset which contains tables has the information regarding various fields related to company hiring data. I have gone through the dataset closely and then executed various queries to extract the required data from the table and also made graphs, charts and tables out of the datasets.

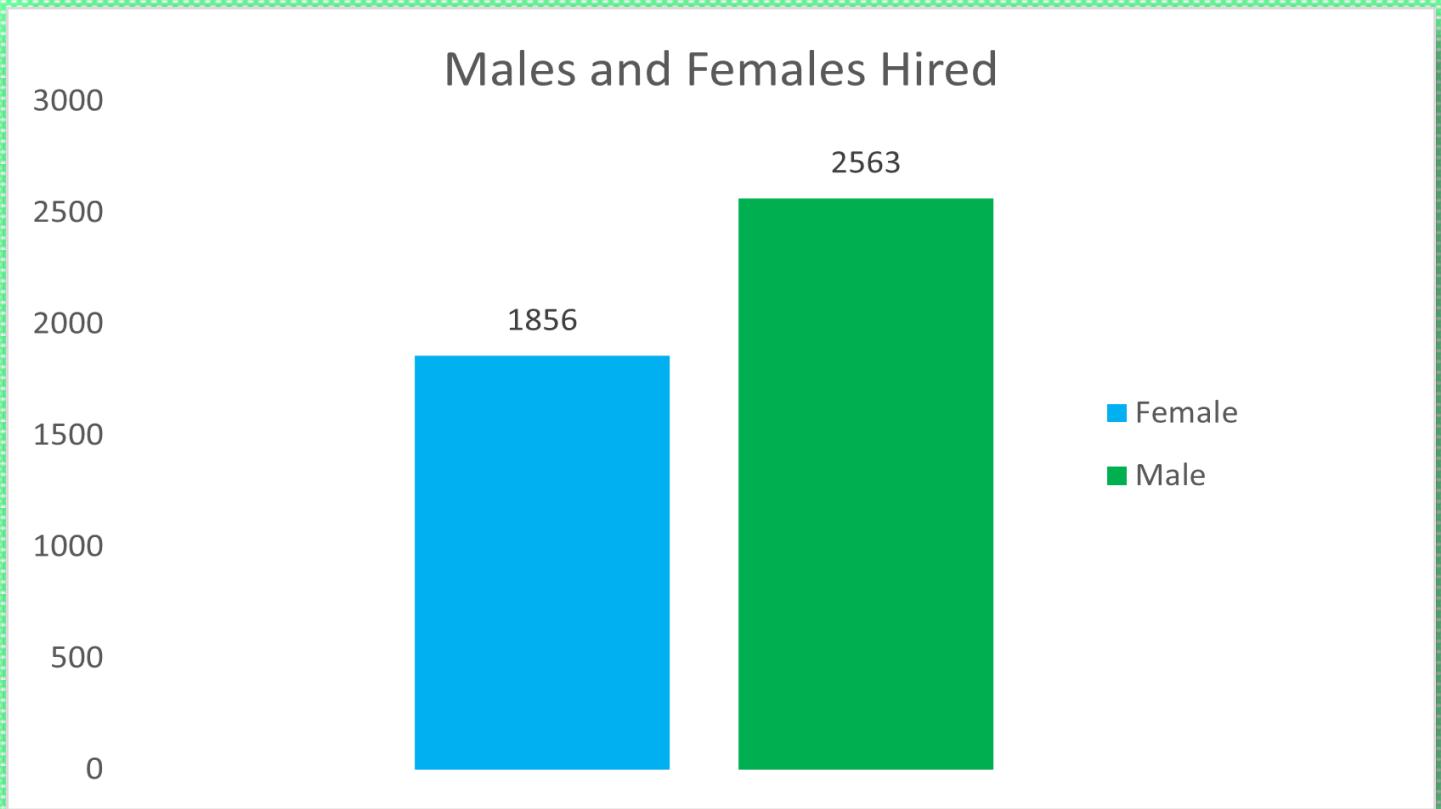
Tech Stack Used:

- PPT – To prepare a detailed report.
- EXCEL – Excel was used to perform entire analysis.

A. Hiring:

How many males and females are Hired?

| Status | Hired | Filter | |
|--------------------------------|--------|---------------|-------------|
| | | Column Labels | Filter |
| Total no male and female hired | Female | Male | Grand Total |
| | 1856 | 2563 | 4419 |



B.Average Salary:

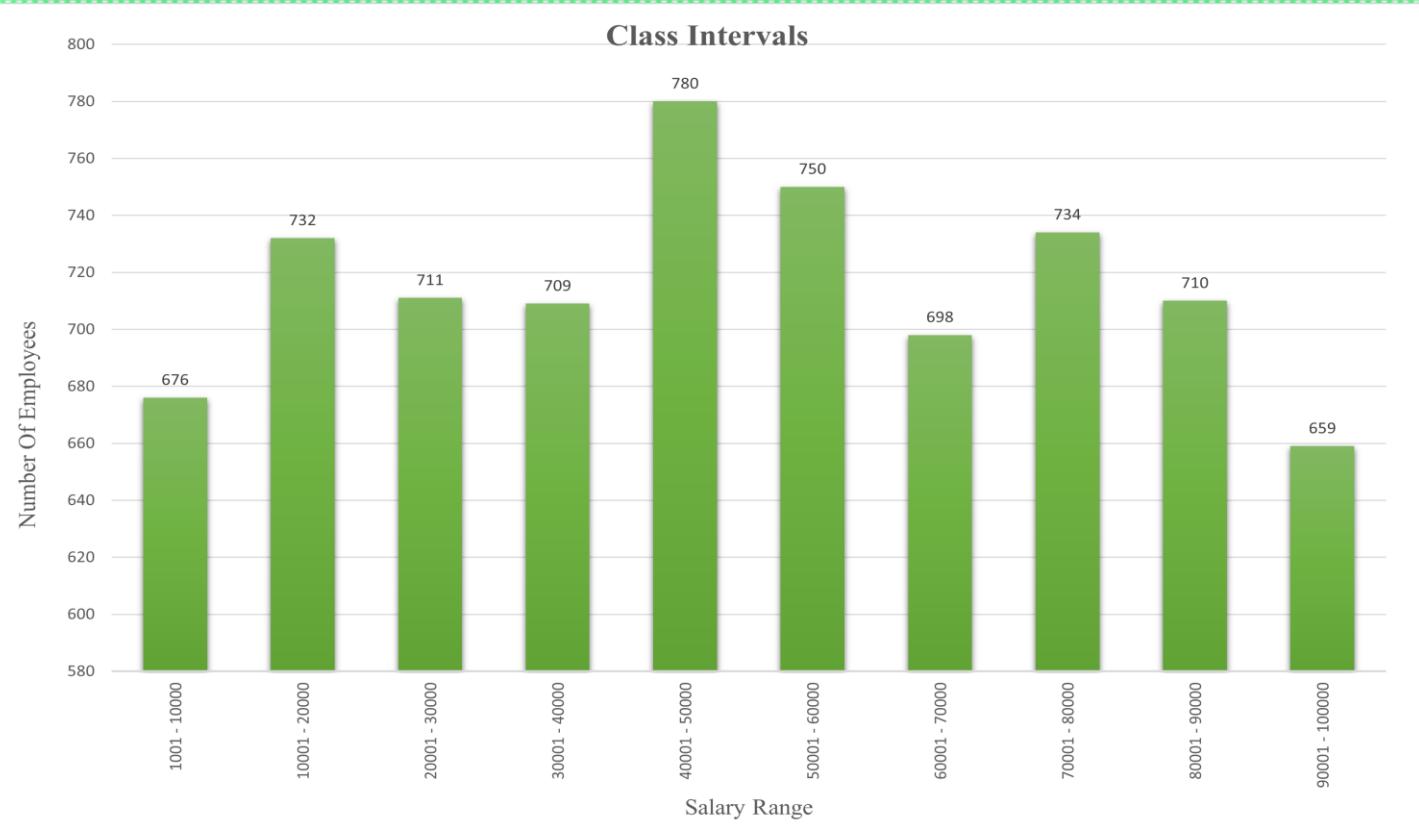
What is the average salary offered in this company?

| Offered Salary | (All) |
|---------------------------|-------|
| Average of Offered Salary | |
| 49983.03 | |

C.Class Intervals:

Draw the class intervals for salary in the company?

| Range | Frequency |
|----------------|-----------|
| 1001 - 10000 | 676 |
| 10001 - 20000 | 732 |
| 20001 - 30000 | 711 |
| 30001 - 40000 | 709 |
| 40001 - 50000 | 780 |
| 50001 - 60000 | 750 |
| 60001 - 70000 | 698 |
| 70001 - 80000 | 734 |
| 80001 - 90000 | 710 |
| 90001 - 100000 | 659 |

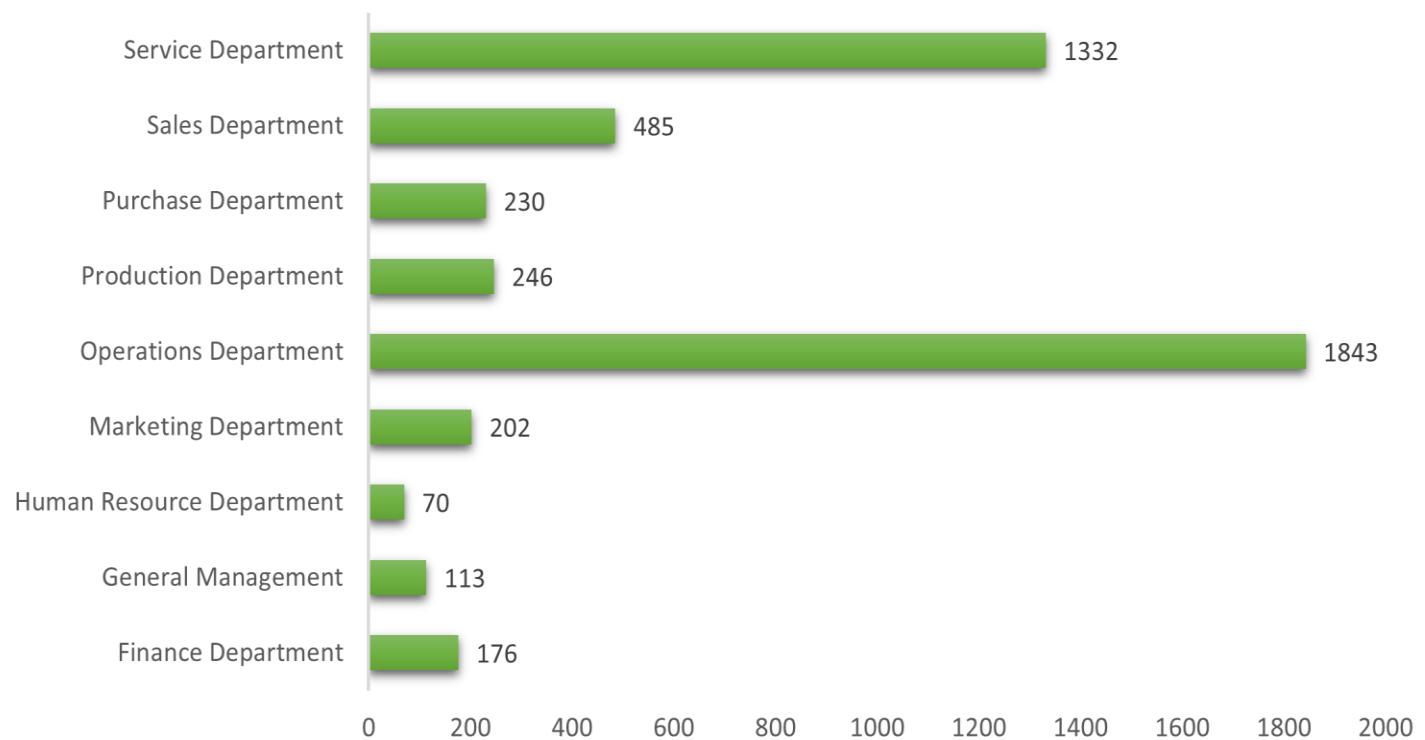


D.Charts and Plots:

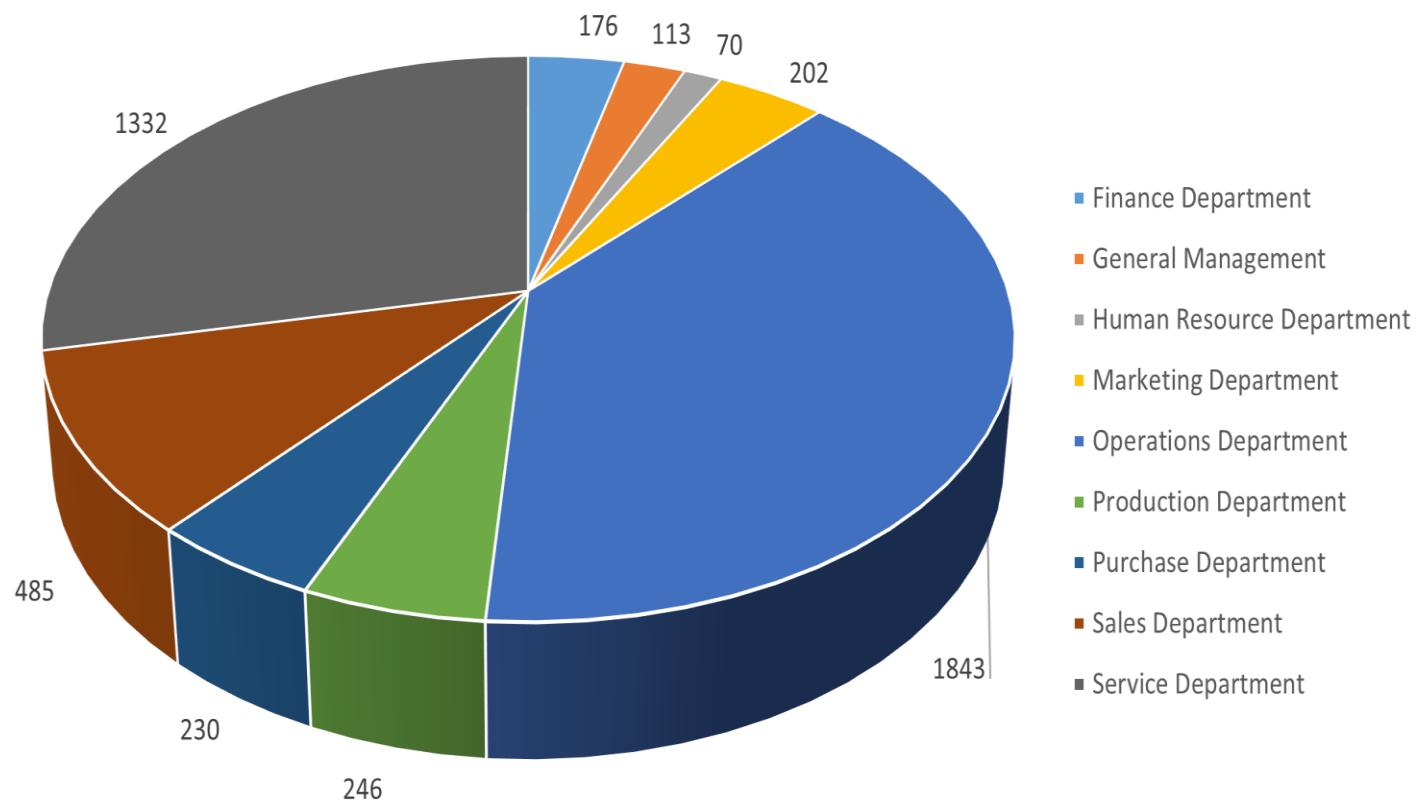
Draw Pie Chart / Bar Graph to show proportion of people working different department?

| Status | Hired | |
|---------------------------|-----------------|--|
| Row Labels | Count of Status | |
| Finance Department | 176 | |
| General Management | 113 | |
| Human Resource Department | 70 | |
| Marketing Department | 202 | |
| Operations Department | 1843 | |
| Production Department | 246 | |
| Purchase Department | 230 | |
| Sales Department | 485 | |
| Service Department | 1332 | |
| Grand Total | 4697 | |

Employees working at different department



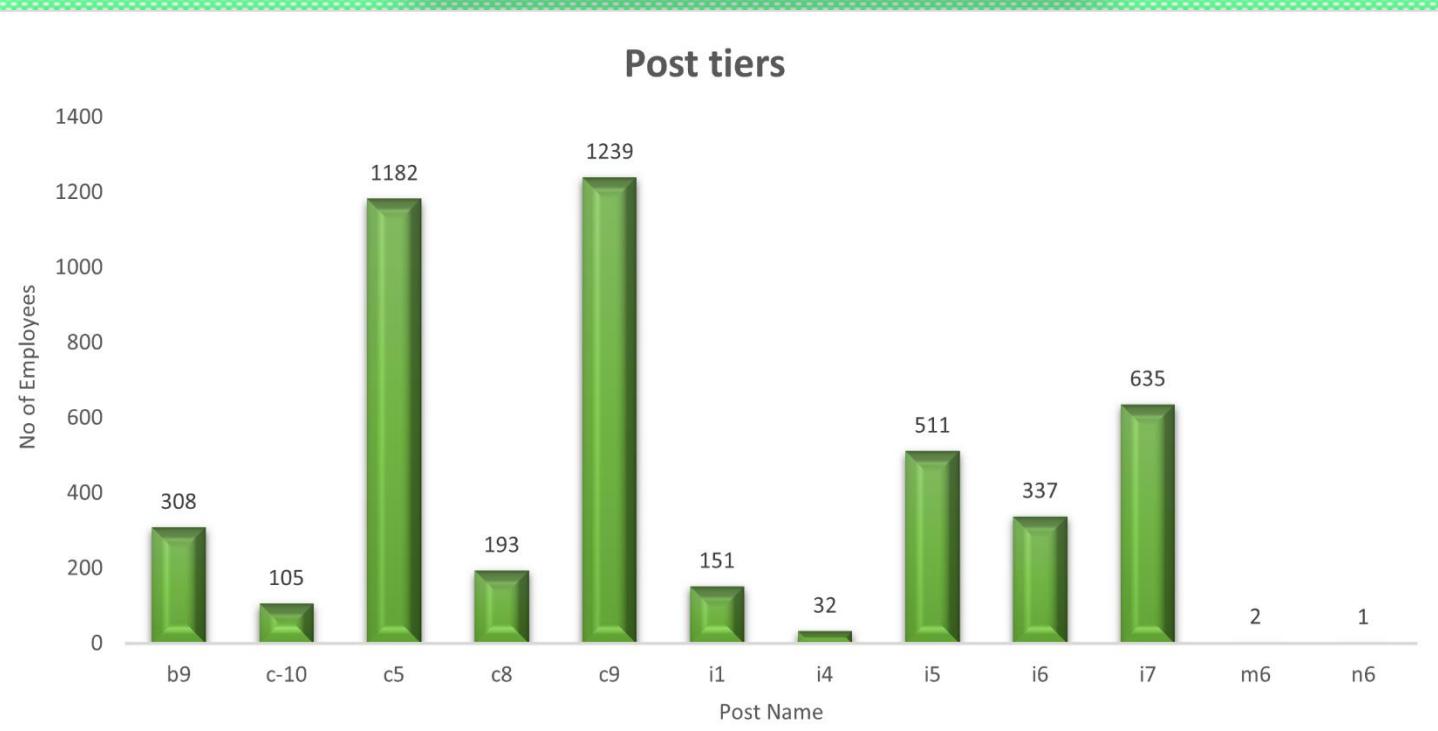
Employees working at different department



E.Charts:

Represent different post tiers using chart/graph?

| Status | Hired |
|--------------------|-------------|
| Row Labels | |
| b9 | 308 |
| c-10 | 105 |
| c5 | 1182 |
| c8 | 193 |
| c9 | 1239 |
| i1 | 151 |
| i4 | 32 |
| i5 | 511 |
| i6 | 337 |
| i7 | 635 |
| m6 | 2 |
| n6 | 1 |
| Grand Total | |
| | 4696 |



Conclusion

Overall, this project provided valuable insights into the hiring process and trends within the company, allowing for improvements to be made to ensure a fair and efficient hiring process.

Hiring Process Analytics helps the company to decide the salaries for new freshers joining the company; also it tells requirement of workforce by each department; it also helps the company decide the appraisals and increment for its current employes.



IMDB Movie Analysis

Project-5



Project Description:

This project is about the IMDB movie analysis based on different criteria. We have IMBD's dataset for different movies, which contains many information and we have to analyze the dataset using the questions that where asked.

Approach:

We have to clean the given datasets which contains tables that contain information regarding various fields related to movies to be rated by IMDB. We have to go through the dataset closely and then executed various queries to extract the required insights from the dataset that was cleaned.

Tech stack used:

PPT – To prepare a detailed report.

EXCEL – Excel was used to perform entire analysis.

Findings:

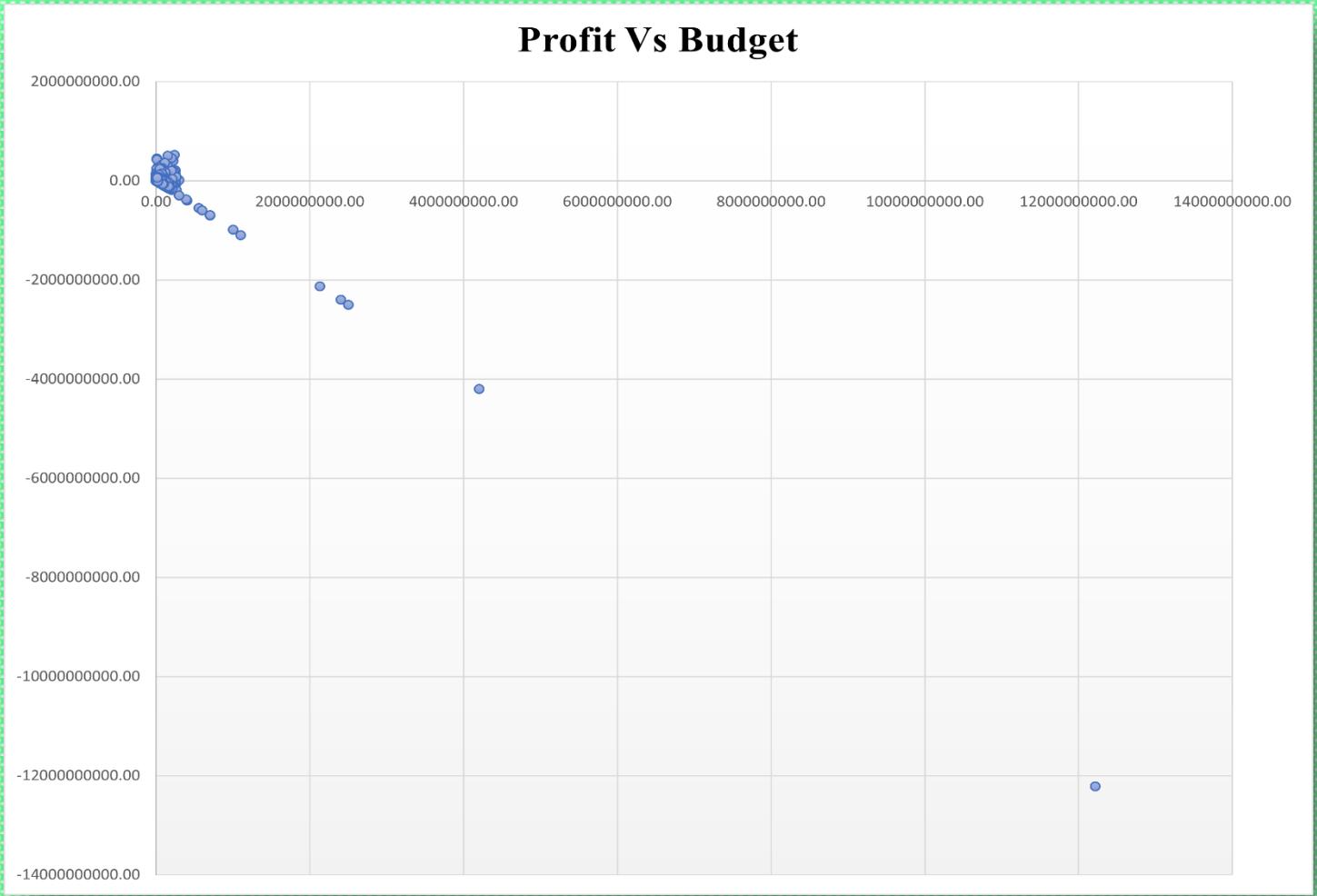
A. Cleaning the data:

1. Dropping unnecessary columns. (Color, director_facebook_likes, actor_3_facebook_likes, actor_2_name, actor_1_facebook_likes, cast_total_facebook_likes, actor_3_name, facenumber_in_posts, plot_keywords, movie_imdb_link, content_rating, actor_2_facebook_likes, aspect_ratio, movie_facebook_likes)

2. Remove Blank Cell / Null Value.

3. Removing Duplicate

B. Movies with highest profit:



OUTLIERS:

- 12213298588
- 4199788333
- 2499804112
- 2397701809

-2127109510

Top 5 Highest Profit

| director_name | actor_1_name | movie_title | title_year | imdb_score | Profit |
|------------------|---------------------|------------------------------------------------|------------|------------|-----------|
| James Cameron | CCH Pounder | Avatar | 2009 | 7.9 | 523505847 |
| Colin Trevorrow | Bryce Dallas Howard | Jurassic World | 2015 | 7 | 502177271 |
| James Cameron | Leonardo DiCaprio | Titanic | 1997 | 7.7 | 458672302 |
| George Lucas | Harrison Ford | Star Wars: Episode V - The Empire Strikes Back | 1977 | 8.7 | 449935665 |
| Steven Spielberg | Henry Thomas | E.T. the Extra-Terrestrial | 1982 | 7.9 | 424449459 |

C. Top 250:

1. Filtering the data where num_voted_users > 25,000.
2. Sort the data using imdb_score column in largest to smallest score.
3. Create column ‘Rank’ for top 250 movies.
4. Use Sequence Formula to Rank.
5. Formula =SEQUENCE(COUNTA(G2:G251),1,1,1)
6. Filter out language by unselecting English to extract the top foreign language movie.

Top 250 Movies

| rank | Profit | imdb_score | title_year | country | budget | language | num_us | num_voted_users | movie_title | actor_1_name | genres | duration | gross | director_name |
|------|-----------|------------|------------|-------------|-----------|------------|--------|-----------------|---------------------------------------------------|---------------------|-----------|----------|-----------|----------------------|
| 1 | 3341469 | 9.3 | 1994 | USA | 25000000 | English | 4144 | 1689764 | The Shawshank Redemption | Morgan Freeman | Crime | 142 | 28341469 | Frank Darabont |
| 2 | 128821952 | 9.2 | 1972 | USA | 6000000 | English | 2238 | 1155770 | The Godfather | Al Pacino | Crime | 175 | 134821952 | Francis Ford Coppola |
| 3 | 348316061 | 9 | 2008 | USA | 185000000 | English | 4667 | 1676169 | The Dark Knight | Christian Bale | Action | 152 | 533316061 | Christopher Nolan |
| 4 | 44300000 | 9 | 1974 | USA | 13000000 | English | 650 | 790926 | The Godfather: Part II | Robert De Niro | Crime | 220 | 57300000 | Francis Ford Coppola |
| 5 | 99930000 | 8.9 | 1994 | USA | 8000000 | English | 2195 | 1324680 | Pulp Fiction | Bruce Willis | Crime | 178 | 107930000 | Quentin Tarantino |
| 6 | 283019252 | 8.9 | 2003 | USA | 94000000 | English | 3189 | 1215718 | The Lord of the Rings: The Return of the King | Orlando Bloom | Action | 192 | 377019252 | Peter Jackson |
| 7 | 74067179 | 8.9 | 1993 | USA | 22000000 | English | 1273 | 865020 | Schindler's List | Liam Neeson | Biography | 185 | 96067179 | Steven Spielberg |
| 8 | 4900000 | 8.9 | 1966 | Italy | 1200000 | Italian | 780 | 503509 | The Good, the Bad and the Ugly | Clint Eastwood | Western | 142 | 6100000 | Sergio Leone |
| 9 | 132568851 | 8.8 | 2010 | USA | 160000000 | English | 2803 | 1468200 | Inception | Leonardo DiCaprio | Action | 148 | 292568851 | Christopher Nolan |
| 10 | -25976605 | 8.8 | 1999 | USA | 63000000 | English | 2968 | 1347461 | Fight Club | Brad Pitt | Drama | 151 | 37023395 | David Fincher |
| 11 | 274691196 | 8.8 | 1994 | USA | 55000000 | English | 1398 | 1251222 | Forrest Gump | Tom Hanks | Comedy | 142 | 329691196 | Robert Zemeckis |
| 12 | 220837577 | 8.8 | 2001 | New Zealand | 93000000 | English | 5060 | 1238746 | The Lord of the Rings: The Fellowship of the Ring | Christopher Lee | Action | 171 | 313837577 | Peter Jackson |
| 13 | 272158751 | 8.8 | 1980 | USA | 18000000 | English | 900 | 837759 | Star Wars: Episode V - The Empire Strikes Back | Harrison Ford | Action | 127 | 290158751 | Irvin Kershner |
| 14 | 108383253 | 8.7 | 1999 | USA | 63000000 | English | 3646 | 1217752 | The Matrix | Keanu Reeves | Action | 136 | 171383253 | Lana Wachowski |
| 15 | 246478898 | 8.7 | 2002 | USA | 94000000 | English | 2417 | 1100446 | The Two Towers | Christopher Lee | Action | 172 | 340478898 | Peter Jackson |
| 16 | 449935665 | 8.7 | 1977 | USA | 11000000 | English | 1470 | 911097 | Star Wars: Episode IV - A New Hope | Harrison Ford | Action | 125 | 460935665 | George Lucas |
| 17 | 21836394 | 8.7 | 1990 | USA | 25000000 | English | 728685 | 7989 | Goodfellas | Robert De Niro | Biography | 146 | 46836394 | Martin Scorsese |
| 18 | 107600000 | 8.7 | 1975 | USA | 4400000 | English | 760 | 680041 | One Flew Over the Cuckoo's Nest | Scatman Crothers | Drama | 133 | 112000000 | Milos Forman |
| 19 | 4263397 | 8.7 | 2002 | Brazil | 3300000 | Portuguese | 749 | 533200 | City of God | Alice Braga | Crime | 135 | 7563397 | Fernando Meirelles |
| 20 | -1730939 | 8.7 | 1954 | Japan | 2000000 | Japanese | 596 | 229012 | Seven Samurai | Takashi Shimura | Action | 202 | 269061 | Akira Kurosawa |
| 21 | 67125340 | 8.6 | 1995 | USA | 33000000 | English | 1080 | 1023511 | Se7en | Morgan Freeman | Crime | 127 | 100125340 | David Fincher |
| 22 | 22991439 | 8.6 | 2014 | USA | 165000000 | English | 2725 | 928227 | Interstellar | Matthew McConaughey | Adventure | 169 | 187991439 | Christopher Nolan |
| 23 | 111727000 | 8.6 | 1991 | USA | 19000000 | English | 916 | 887467 | The Silence of the Lambs | Anthony Hopkins | Crime | 138 | 130727000 | Jonathan Demme |
| 24 | 146119491 | 8.6 | 1998 | USA | 70000000 | English | 2277 | 881236 | Saving Private Ryan | Tom Hanks | Action | 169 | 216119491 | Steven Spielberg |
| 25 | -787759 | 8.6 | 1998 | USA | 7500000 | English | 1420 | 782437 | American History X | Ethan Suplee | Crime | 101 | 6712241 | Tony Kaye |
| 26 | 17272306 | 8.6 | 1995 | USA | 6000000 | English | 1182 | 740918 | The Usual Suspects | Kevin Spacey | Crime | 106 | 23272306 | Bryan Singer |
| 27 | 8950114 | 8.6 | 2001 | Japan | 10000000 | Japanese | 102 | 447074 | Grave of the Fireflies | Mitsuru Fukikoshi | Animation | 125 | 10049886 | Hayao Miyazaki |

Top Foreign Language Movies

| rank | Profit | imdb_score | title_year | country | budget | language | num_use | num_voted | movie_title | actor_1_name | genres | duration | gross | director_name |
|------|-------------|------------|------------|--------------|----------|------------|---------|-----------|--------------------------------|----------------------|-----------|----------|----------|----------------------------------|
| 1 | 4900000 | 8.9 | 1966 | Italy | 1200000 | Italian | 780 | 503509 | The Good, the Bad and the Ugly | Clint Eastwood | Western | 142 | 6100000 | Sergio Leone |
| 2 | 4263397 | 8.7 | 2002 | Brazil | 3300000 | Portuguese | 749 | 533200 | City of God | Alice Braga | Crime | 135 | 7563397 | Fernando Meirelles |
| 3 | -1730939 | 8.7 | 1954 | Japan | 2000000 | Japanese | 596 | 229012 | Seven Samurai | Takashi Shimura | Action | 202 | 269061 | Akira Kurosawa |
| 4 | -8950114 | 8.6 | 2001 | Japan | 19000000 | Japanese | 902 | 417971 | Spirited Away | Bunta Sugawara | Animation | 125 | 10049886 | Hayao Miyazaki |
| 5 | 9284657 | 8.5 | 2006 | Germany | 2000000 | German | 407 | 259379 | The Lives of Others | Sebastian Koch | Thriller | 137 | 11284657 | Florian Henckel von Donnersmarck |
| 6 | 745402 | 8.5 | 1997 | Iran | 180000 | Persian | 130 | 27882 | Children of Heaven | Bahare Seddiqi | Drama | 89 | 925402 | Majid Majidi |
| 7 | -43798339 | 8.4 | 2001 | France | 77000000 | French | 1314 | 534262 | Amélie | Mathieu Kassovitz | Romance | 122 | 33201661 | Jean-Pierre Jeunet |
| 8 | -818710 | 8.4 | 2003 | South Korea | 3000000 | Korean | 809 | 356181 | Oldboy | Min-sik Choi | Drama | 120 | 2181290 | Chan-wook Park |
| 9 | -2397701809 | 8.4 | 1997 | Japan | 2.4E+09 | Japanese | 570 | 221552 | Princess Mononoke | Miyoko Takeuchi | Animation | 134 | 2298191 | Hayao Miyazaki |
| 10 | -2566866 | 8.4 | 1981 | West Germany | 14000000 | German | 426 | 168203 | Das Boot | Jürgen Prochnow | Adventure | 293 | 11433134 | Wolfgang Petersen |
| 11 | 6598492 | 8.4 | 2011 | Iran | 500000 | Persian | 264 | 151812 | A Separation | Shahab Hosseini | Drama | 123 | 7098492 | Asghar Farhadi |
| 12 | -11528148 | 8.4 | 2015 | India | 18026148 | Telugu | 410 | 62756 | Baahubali: The Beginning | Prabhas | Action | 159 | 6498000 | S.S. Rajamouli |
| 13 | -7998060 | 8.3 | 2004 | Germany | 13500000 | German | 564 | 248354 | Downfall | Thomas Kretschmann | Biography | 178 | 5501940 | Oliver Hirschbiegel |
| 14 | -3189032 | 8.3 | 2012 | Denmark | 3800000 | Danish | 249 | 170155 | The Hunt | Thomas Bo Larsen | Drama | 115 | 610968 | Thomas Vinterberg |
| 15 | -5973565 | 8.3 | 1927 | Germany | 6000000 | German | 413 | 111841 | Metropolis | Brigitte Helm | Sci-Fi | 145 | 26435 | Fritz Lang |
| 16 | 24123143 | 8.2 | 2006 | Spain | 13500000 | Spanish | 1083 | 467234 | Pan's Labyrinth | Ivana Baquero | Fantasy | 112 | 37623143 | Guillermo del Toro |
| 17 | -19289545 | 8.2 | 2004 | Japan | 24000000 | Japanese | 330 | 214091 | Howl's Moving Castle | Hiromasa Yonebayashi | Animation | 119 | 4710455 | Hayao Miyazaki |
| 18 | -3189032 | 8.2 | 2009 | Argentina | 2000000 | Spanish | 231 | 131831 | The Secret in Their Eyes | Ricardo Darín | Mystery | 129 | 20167424 | Juan José Campanella |
| 19 | 57096 | 8.2 | 2010 | Canada | 6800000 | French | 156 | 80429 | Incendies | Slava Polonsky | Mystery | 139 | 6857096 | Denis Villeneuve |
| 20 | 3383834 | 8.1 | 2000 | Mexico | 2000000 | Spanish | 361 | 173551 | Amores Perros | Adriana Barraza | Thriller | 115 | 5383834 | Alejandro G. Iñárrizcoa |
| 21 | -1099560838 | 8.1 | 1988 | Japan | 1.1E+09 | Japanese | 430 | 106160 | Akira | Mitsuo Iwata | Animation | 124 | 439162 | Katsuhiro Ōtomo |
| 22 | -3991940 | 8.1 | 2007 | Brazil | 4000000 | Portuguese | 107 | 81644 | Elite Squad | Wagner Moura | Crime | 115 | 8060 | José Padilha |
| 23 | 347780 | 8.1 | 1998 | Denmark | 1300000 | Danish | 258 | 65951 | The Celebration | Ulrich Thomsen | Drama | 105 | 1647780 | Thomas Vinterberg |

D. Best Directors:

| Top 10 Directors | | Average of imdb_score |
|-------------------|--|-----------------------|
| Charles Chaplin | | 8.6 |
| Tony Kaye | | 8.6 |
| Alfred Hitchcock | | 8.5 |
| Damien Chazelle | | 8.5 |
| Majid Majidi | | 8.5 |
| Ron Fricke | | 8.5 |
| Sergio Leone | | 8.4 |
| Christopher Nolan | | 8.4 |
| Asghar Farhadi | | 8.4 |
| Marius A. Markevi | | 8.4 |

E. Popular Genres:

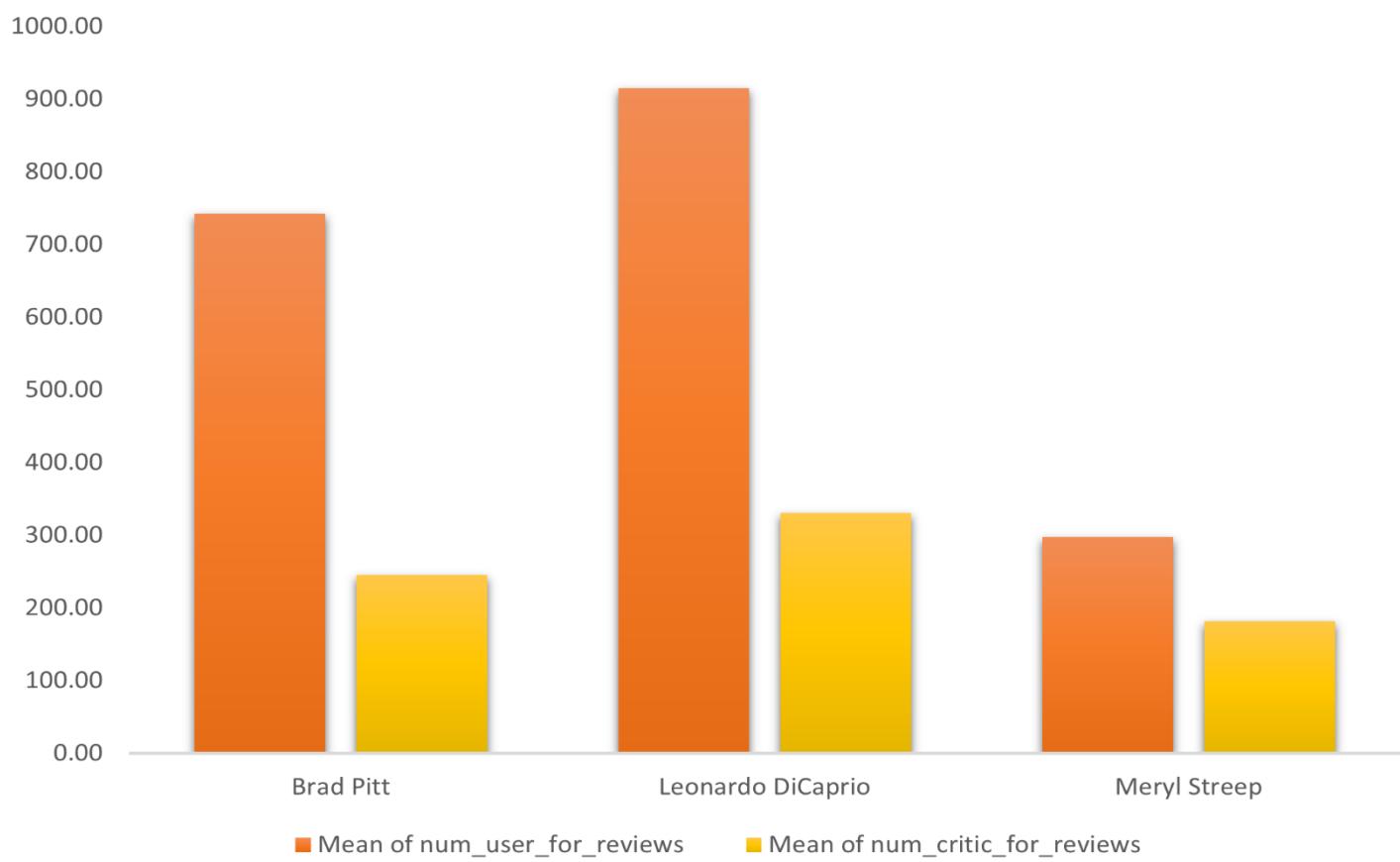
| Popular Genres | Average of imdb_score |
|-------------------------------------------------|-----------------------|
| Crime Drama Fantasy Mystery | 8.5 |
| Adventure Animation Drama Family Musical | 8.5 |
| Adventure Drama Thriller War | 8.4 |
| Adventure Animation Fantasy | 8.4 |
| Action Adventure Drama Fantasy War | 8.4 |
| Documentary War | 8.3 |
| Documentary Drama Sport | 8.3 |
| Biography Drama History Music | 8.3 |
| Adventure Animation Comedy Drama Family Fantasy | 8.3 |
| Adventure Drama War | 8.3 |

F. Charts:

1. Highest Mean

| actor_1_name | Mean of num_user_for_reviews | Mean of num_critic_for_reviews |
|-------------------|------------------------------|--------------------------------|
| Brad Pitt | 742.35 | 245 |
| Leonardo DiCaprio | 914.48 | 330.19 |
| Meryl Streep | 297.18 | 181.45 |

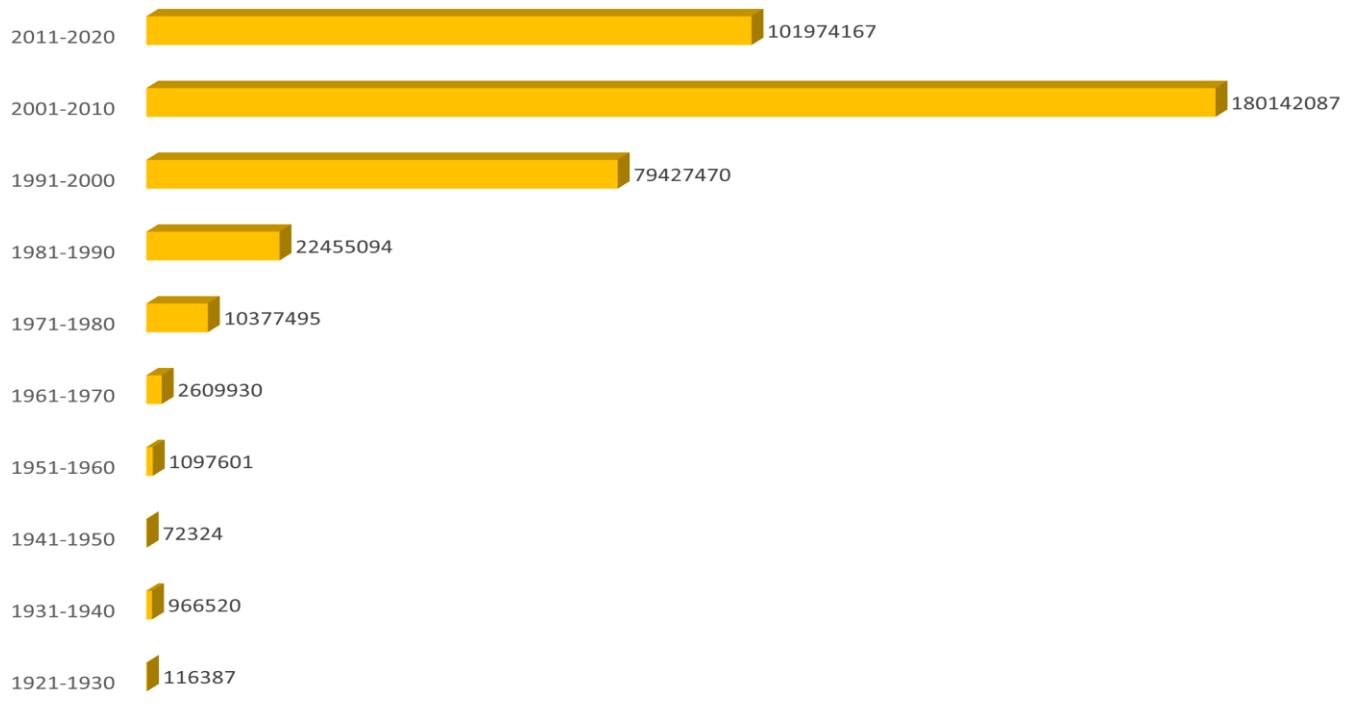
Highest Mean



2. Voted users over decades

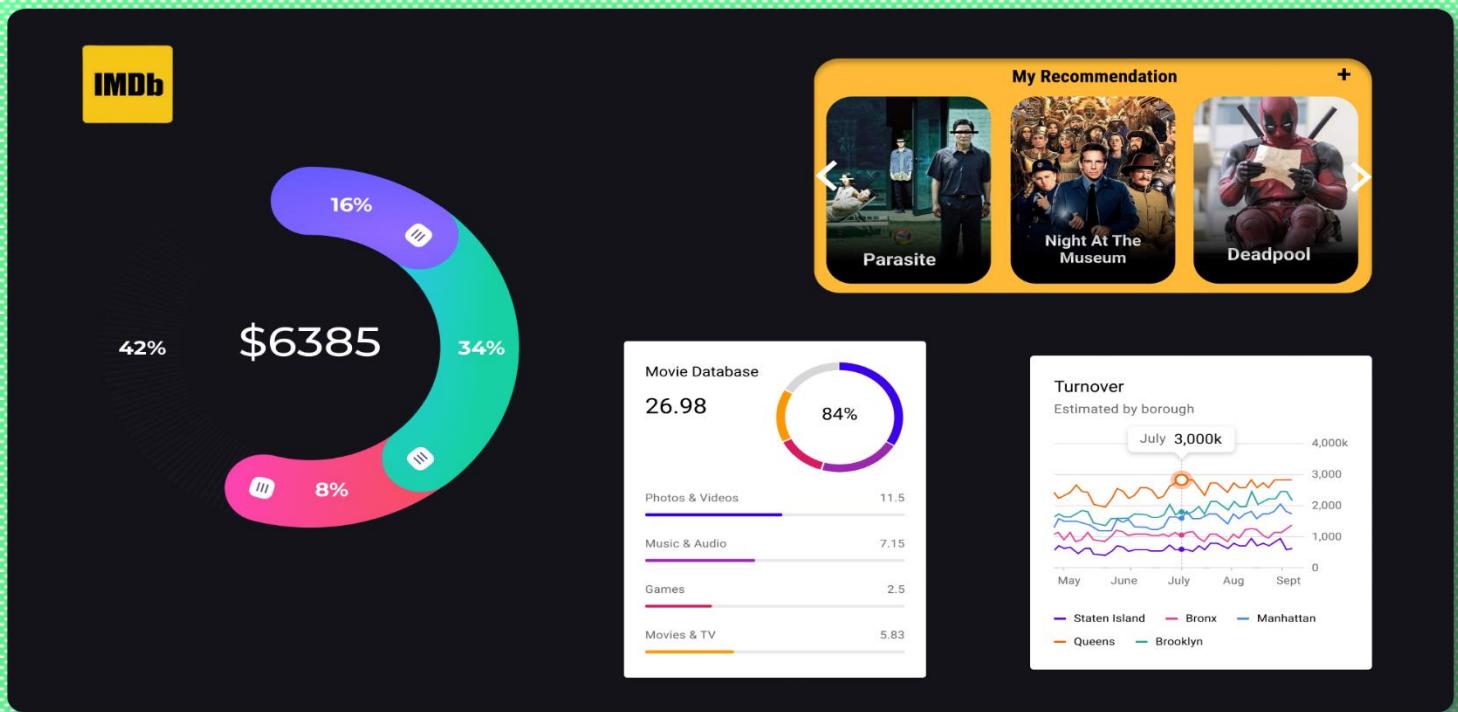
| Decade | Number of Users Voted |
|-----------|-----------------------|
| 1921-1930 | 116387 |
| 1931-1940 | 966520 |
| 1941-1950 | 72324 |
| 1951-1960 | 1097601 |
| 1961-1970 | 2609930 |
| 1971-1980 | 10377495 |
| 1981-1990 | 22455094 |
| 1991-2000 | 79427470 |
| 2001-2010 | 180142087 |
| 2011-2020 | 101974167 |

Number of Users Voted over Decades



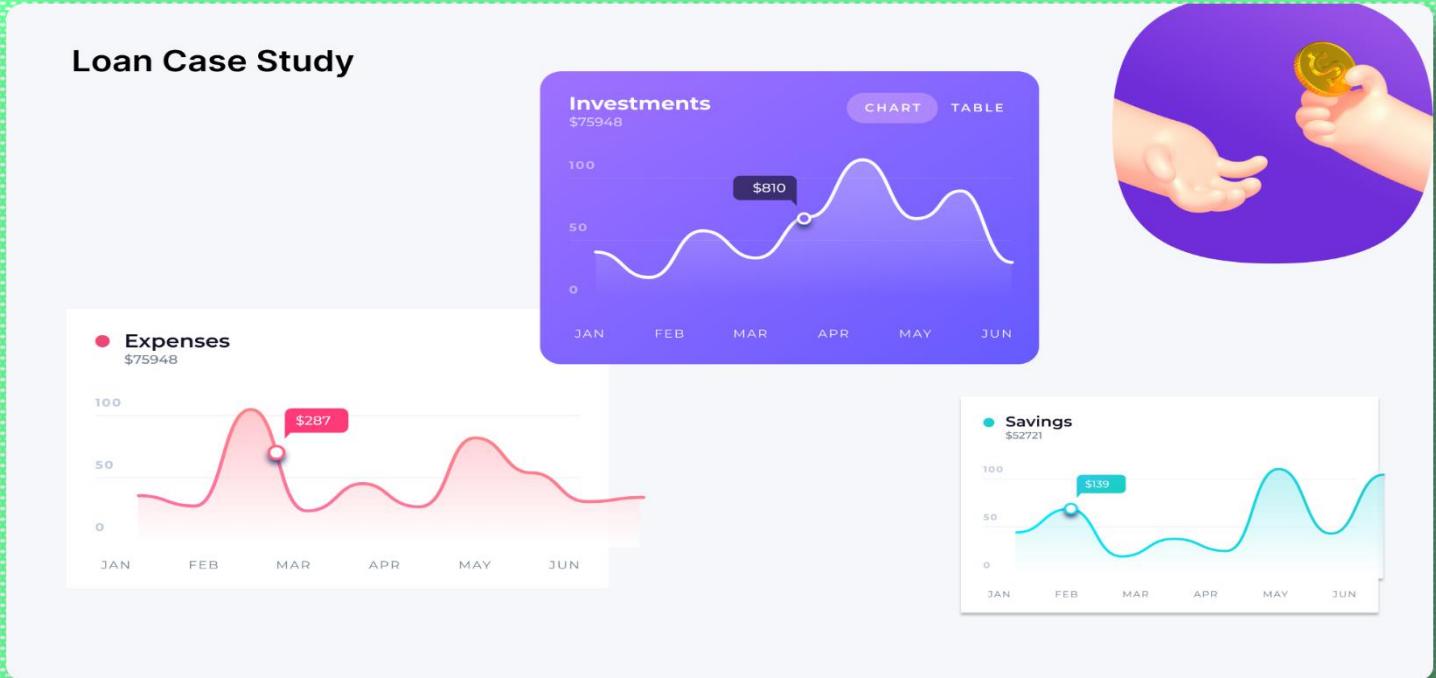
Conclusion

Overall, this project provided valuable insights into IMDB Movie Analysis such analysis is done not only by Movie makers before movie production, but it is also done by various investors, stakeholders, theatre outlet owners, allowing for improvements to be made to ensure a fair and efficient Profits.



Bank Loan Case Study

Project-6



Project Description:

This project is about a case study on Bank Loan. We have to follow EDA (Exploratory Data Analysis) based on our analysis, we will get the insights for required questions that were raised.

Business Objectives:

The main aim of this project is to identify patterns that indicate if a customer will have difficulty paying their installments. This information can be used to make

decisions such as denying the loan, reducing the amount of loan, or lending at a higher interest rate to risky applicants. The company wants to understand the key factors behind loan default so it can make better decisions about loan approval.

Approach:

I have imported the dataset to perform analysis after that we have to handle the data by clean, replace inappropriate values and preprocess the data, then after modification of dataset we had to perform necessary methods to get the desired output/result.

Tech/Tools used:

PPT – To prepare a detailed report.

EXCEL – Excel was used to perform entire analysis.

Insights:

A. Identify Missing Data and Deal with it Appropriately:

1. The dataset have been cleaned and analyzed in the excel,
2. To clean the data the percentage of null values needs to be analyzed and those columns that have more than 50% of the null data have to be dropped,

3. And those columns with less than 50% of the null data have to be replaced with mean or median or the highest occurring variables,

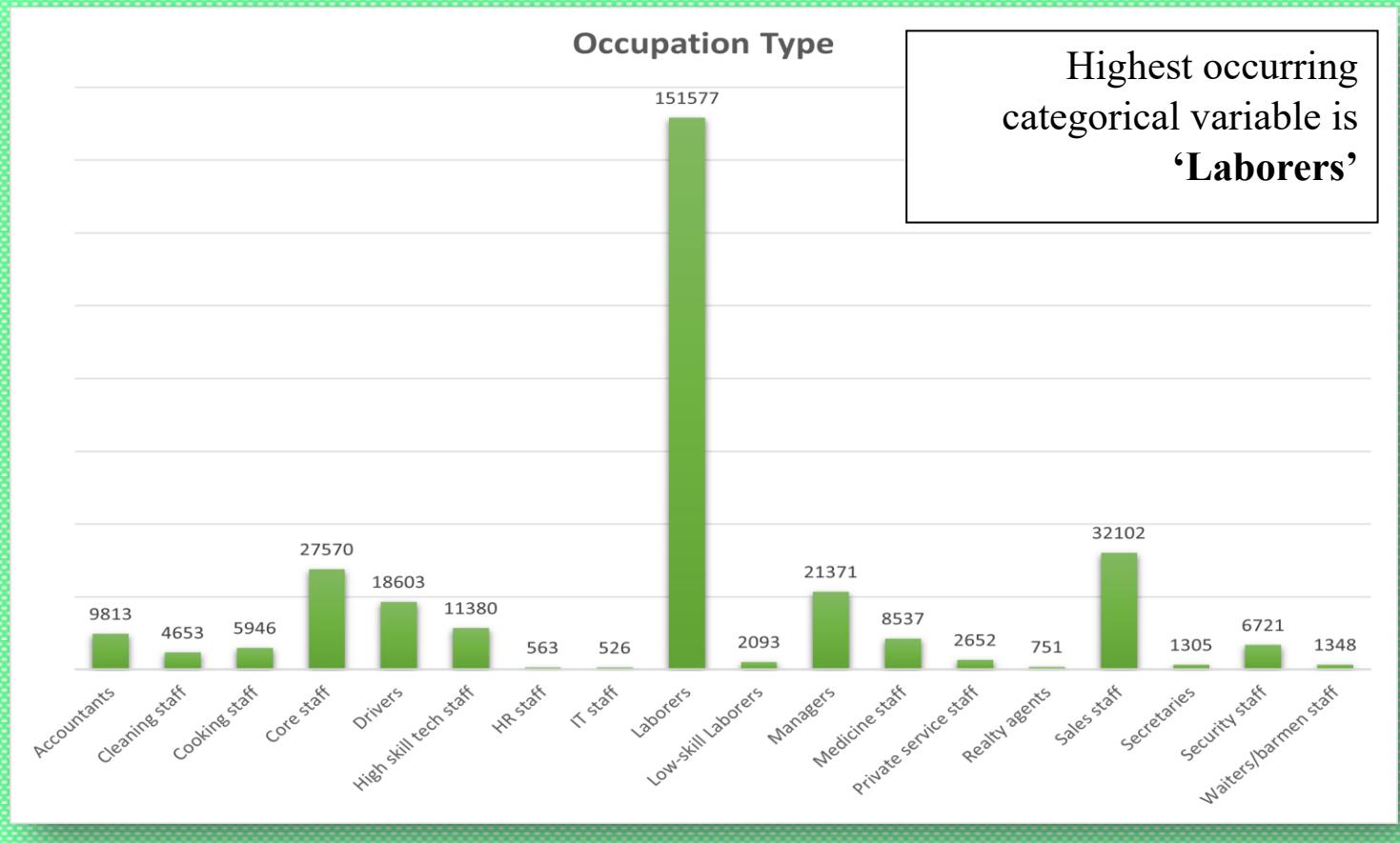
4. 50% is the threshold that I have chosen.

ALL THE COLUMN NAME WHICH ARE IN THE FOLLOWING TABLE ARE NEED TO BE DROPPED DOWN AS THEY HAVE NULL VALUES GREATER THAN OR EQUAL TO 50% AND THE COLUMN THAT ARE IRRELEVANT COLUMNS FOR DOING OUR ANALYSIS.

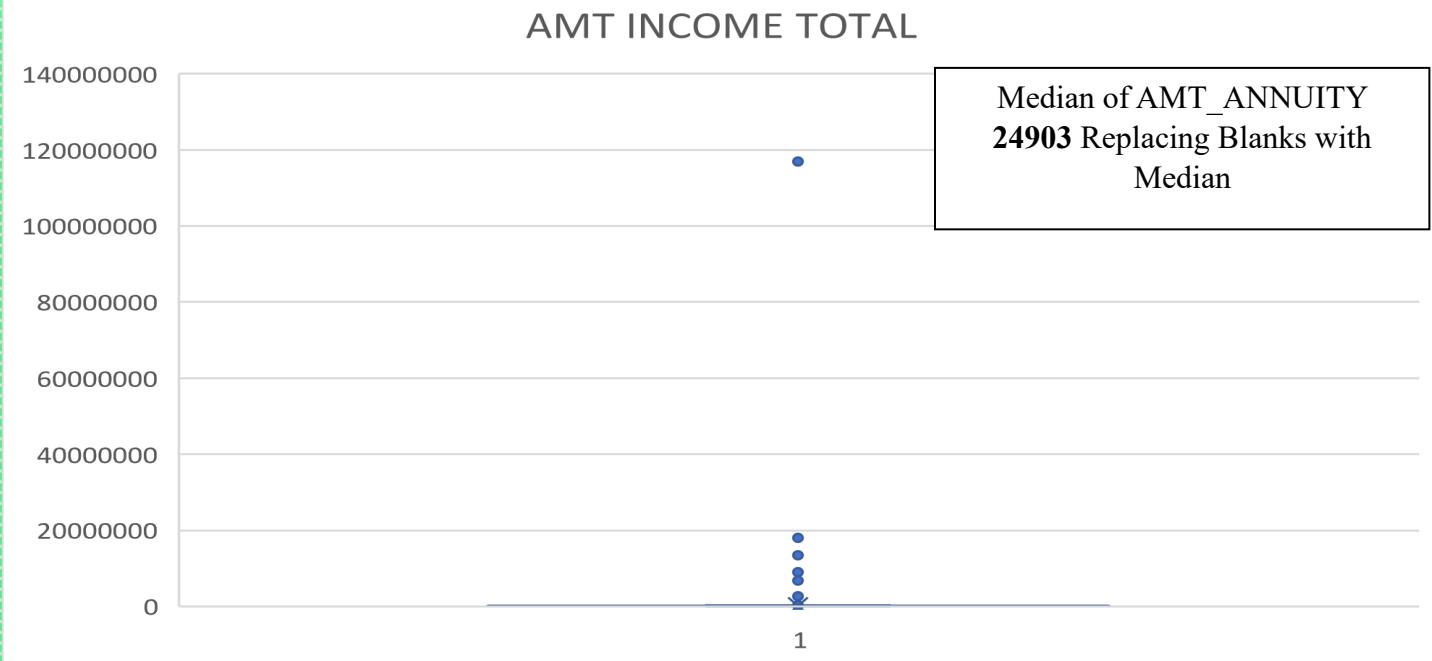
| Column name | Total number of null values | Percentage of null value in that column | ROUND PER |
|----------------------------|-----------------------------|-----------------------------------------|-----------|
| OWN_CAR_AGE | 202930 | 65.99113528 | 66 |
| EXT_SOURCE_1 | 173379 | 56.38139774 | 56 |
| APARTMENTS_AVG | 156061 | 50.74972928 | 51 |
| BASEMENTAREA_AVG | 179943 | 58.51595553 | 59 |
| YEARS_BUILD_AVG | 204488 | 66.49778382 | 66 |
| COMMON_AREA_AVG | 214865 | 69.87229725 | 70 |
| ELEVATORS_AVG | 163891 | 53.29597966 | 53 |
| ENTRANCES_AVG | 154828 | 49.70488861 | 50 |
| FLOORSMAX_AVG | 153021 | 49.76114676 | 50 |
| FLOORSMIN_AVG | 208642 | 67.84862981 | 68 |
| LANDAREA_AVG | 182590 | 59.37673774 | 59 |
| LIVINGAPARTMENTS_AVG | 210199 | 68.35495316 | 68 |
| LIVINGAREA_AVG | 154350 | 50.19332642 | 50 |
| NONLIVINGAPARTMENTS_AVG | 213514 | 69.43296337 | 69 |
| NONLIVINGAREA_AVG | 169682 | 55.17916432 | 55 |
| APARTMENTS_MODE | 156061 | 50.74972928 | 51 |
| BASEMENTAREA_MODE | 179943 | 58.51595553 | 59 |
| YEARS_BUILD_MODE | 204488 | 66.49778382 | 66 |
| COMMON_AREA_MODE | 214865 | 69.87229725 | 70 |
| ELEVATORS_MODE | 163891 | 53.29597966 | 53 |
| ENTRANCES_MODE | 154828 | 50.34876801 | 50 |
| FLOORSMAX_MODE | 153020 | 49.76082156 | 50 |
| FLOORSMIN_MODE | 208642 | 67.84862981 | 68 |
| LANDAREA_MODE | 182590 | 59.37673774 | 59 |
| LIVINGAPARTMENTS_MODE | 210199 | 68.35495316 | 68 |
| LIVINGAREA_MODE | 154350 | 50.19332642 | 50 |
| NONLIVINGAPARTMENTS_MODE | 213514 | 69.43296337 | 69 |
| NONLIVINGAREA_MODE | 169682 | 55.17916432 | 55 |
| APARTMENTS_MEDIAN | 156061 | 50.74972928 | 51 |
| BASEMENTAREA_MEDIAN | 179943 | 58.51595553 | 59 |
| YEARS_BUILD_MEDIAN | 204488 | 66.49778382 | 66 |
| COMMON_AREA_MEDIAN | 214865 | 69.87229725 | 70 |
| ELEVATORS_MEDIAN | 163891 | 53.29597966 | 53 |
| ENTRANCES_MEDIAN | 154828 | 50.34876801 | 50 |
| FLOORSMAX_MEDIAN | 153020 | 49.76082156 | 50 |
| FLOORSMIN_MEDIAN | 208642 | 67.84862981 | 68 |
| LANDAREA_MEDIAN | 182590 | 59.37673774 | 59 |
| LIVINGAPARTMENTS_MEDIAN | 210199 | 68.35495316 | 68 |
| LIVINGAREA_MEDIAN | 154350 | 50.19332642 | 50 |
| NONLIVINGAPARTMENTS_MEDIAN | 213514 | 69.43296337 | 69 |
| NONLIVINGAREA_MEDIAN | 169682 | 55.17916432 | 55 |
| FONDKAPREMONT_MODE | 210295 | 68.38617155 | 68 |
| HOUSETYPE_MODE | 154297 | 50.17609126 | 50 |
| WALLSMATERIAL_MODE | 156341 | 50.84078293 | 5 |

| Column name | Total number of null values | Percentage of null value | ROUND PER |
|-------------------------------|-----------------------------|--------------------------|-----------|
| FLAG_MOBIL | 1 | 0.000325192 | 0 |
| FLAG_EMPLOY_PHONE | 55387 | 18.01138821 | 18 |
| FLAG_WORK_PHONE | 0 | 0 | 0 |
| FLAG_CONT_MOBILE | 0 | 0 | 0 |
| FLAG_PHONE | 0 | 0 | 0 |
| FLAG_EMAIL | 0 | 0 | 0 |
| CNT_FAMILY_MEMBERS | 2 | 0.000650383 | 0 |
| REGION_RATING_CLIENT | 0 | 0 | 0 |
| REGION_RATING_CLIENT_W_CITY | 0 | 0 | 0 |
| EXT_SOURCE_3 | 60965 | 19.82530706 | 20 |
| YEAR_BEGINEXPLUATATION_AVG | 150008 | 48.78134441 | 49 |
| YEAR_BEGINEXPLUATATION_MODE | 150007 | 48.78101922 | 49 |
| YEAR_BEGINEXPLUATATION_MEDIAN | 150007 | 48.78101922 | 49 |
| TOTAL_AREA_MODE | 148431 | 48.26851722 | 48 |
| EMERGENCYSTATE_MODE | 145755 | 47.39830445 | 47 |
| DAYS_LAST_PHONE_CHANGE | 1 | 0.000325192 | 0 |
| FLAG_DOC_2 | 0 | 0 | 0 |
| FLAG_DOC_3 | 0 | 0 | 0 |
| FLAG_DOC_4 | 0 | 0 | 0 |
| FLAG_DOC_5 | 0 | 0 | 0 |
| FLAG_DOC_6 | 0 | 0 | 0 |
| FLAG_DOC_7 | 0 | 0 | 0 |
| FLAG_DOC_8 | 0 | 0 | 0 |
| FLAG_DOC_9 | 0 | 0 | 0 |
| FLAG_DOC_10 | 0 | 0 | 0 |
| FLAG_DOC_11 | 0 | 0 | 0 |
| FLAG_DOC_12 | 0 | 0 | 0 |
| FLAG_DOC_13 | 0 | 0 | 0 |
| FLAG_DOC_14 | 0 | 0 | 0 |
| FLAG_DOC_15 | 0 | 0 | 0 |
| FLAG_DOC_16 | 0 | 0 | 0 |
| FLAG_DOC_17 | 0 | 0 | 0 |
| FLAG_DOC_18 | 0 | 0 | 0 |
| FLAG_DOC_19 | 0 | 0 | 0 |
| FLAG_DOC_20 | 0 | 0 | 0 |
| FLAG_DOC_21 | 0 | 0 | 0 |

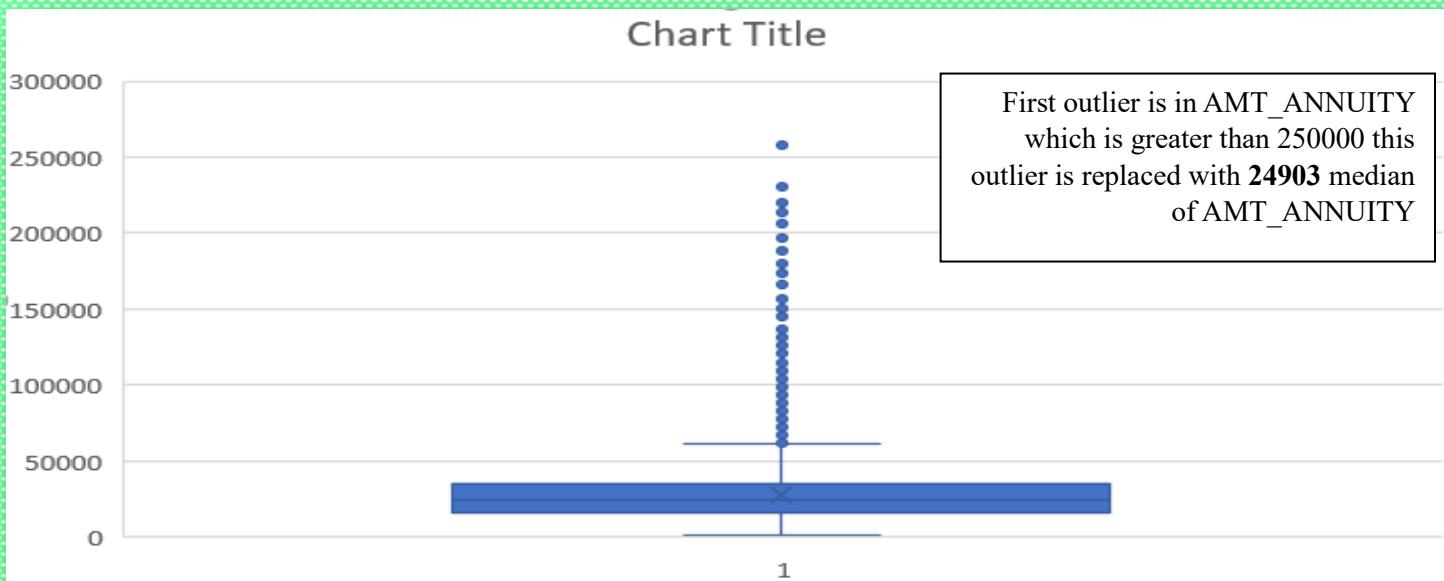
Replacing Blank cell in Occupation_Type column of the Application Dataset with the highest occurring variable.



Replacing Blanks in AMT_ANNUITY column of the Application Dataset with the median of the AMT_ANNUITY as there exists outliers in the AMT_ANNUITY column.



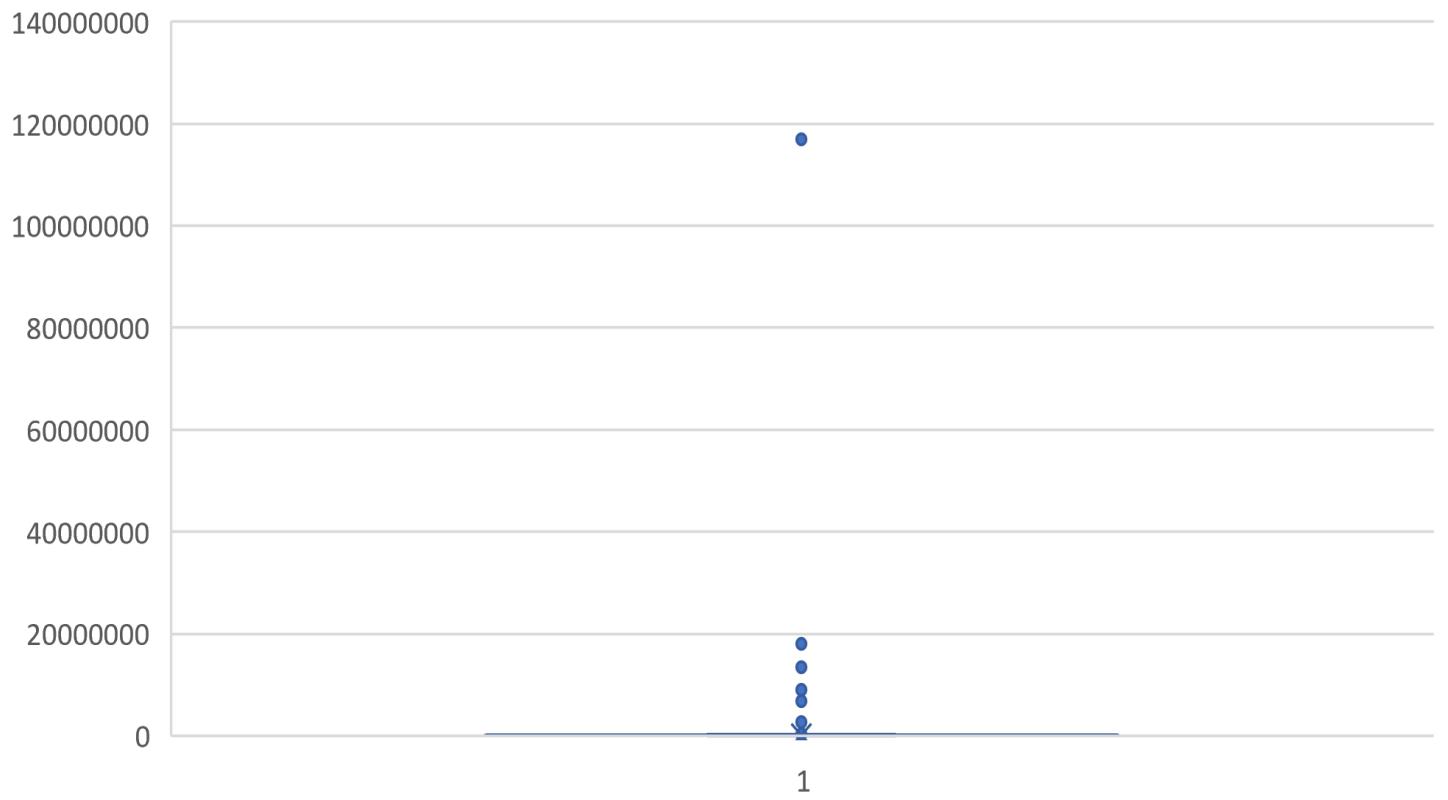
B. Identify Outliers in the Dataset:



Here we can observe that there is huge difference between the 25%, 50% and 75% quartile and this is due to presence of outliers But since the amount of total income varies from person to person we will not remove the outlier.

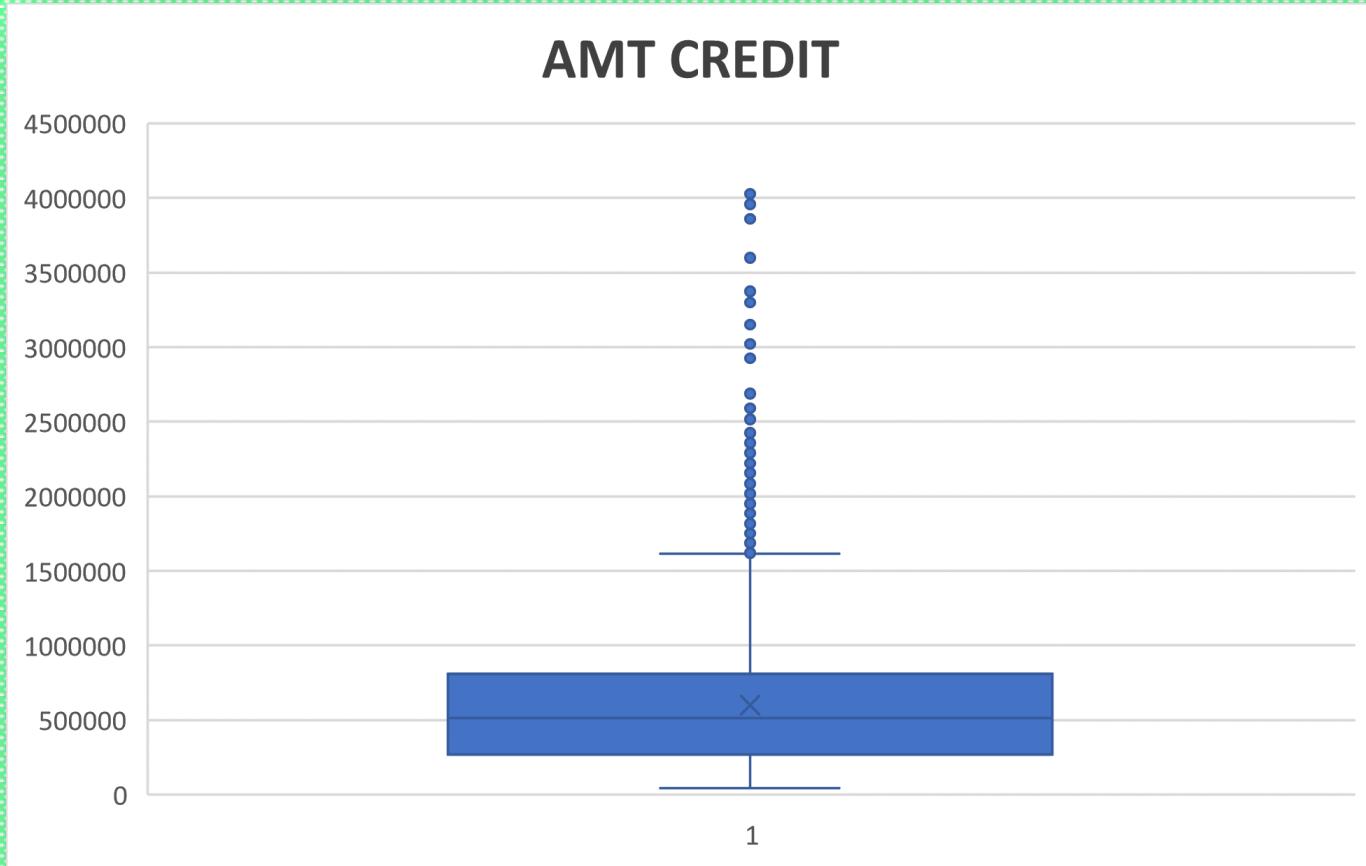
| | Quartiles at AMT_INCOME_TOTAL |
|-----|-------------------------------|
| MIN | 25650 |
| 25% | 112500 |
| 50% | 147150 |
| 75% | 202500 |
| MAX | 117000000 |

AMT INCOME TOTAL



| AMT_CREDIT | |
|------------|-------------------------|
| | Quartiles at AMT_CREDIT |
| MIN | 45000 |
| 25% | 270000 |
| 50% | 513531 |
| 75% | 808650 |
| MAX | 4050000 |

From the chart it is clear that outliers lie in the 98% and near the max side of box plot and also there is a significant difference between the 75% quartile and the max value and this is due the presence of the outliers, But since the amount of credit varies from person to person hence we will not remove the outliers.

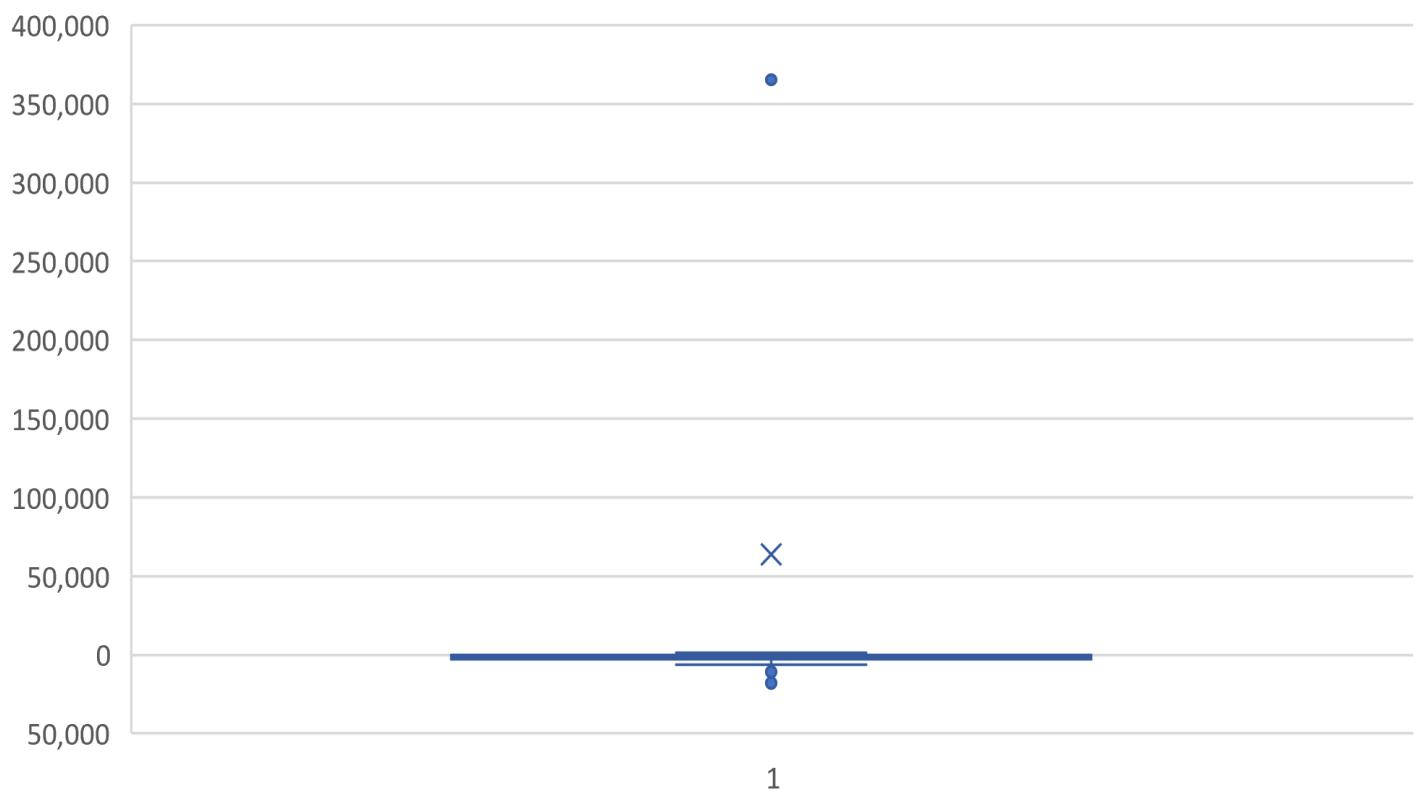


| | DAYS_EMPLOYED |
|----------------------------|---------------|
| Quartiles at DAYS_EMPLOYED | |
| MAX | 17912.00 |
| 75% | 2760.00 |
| 50% | 1213.00 |
| 25% | 289.00 |
| MIN | 365243.00 |

There exists only 1 outlier i.e. + or - 365243

Replace with median **1213.00**

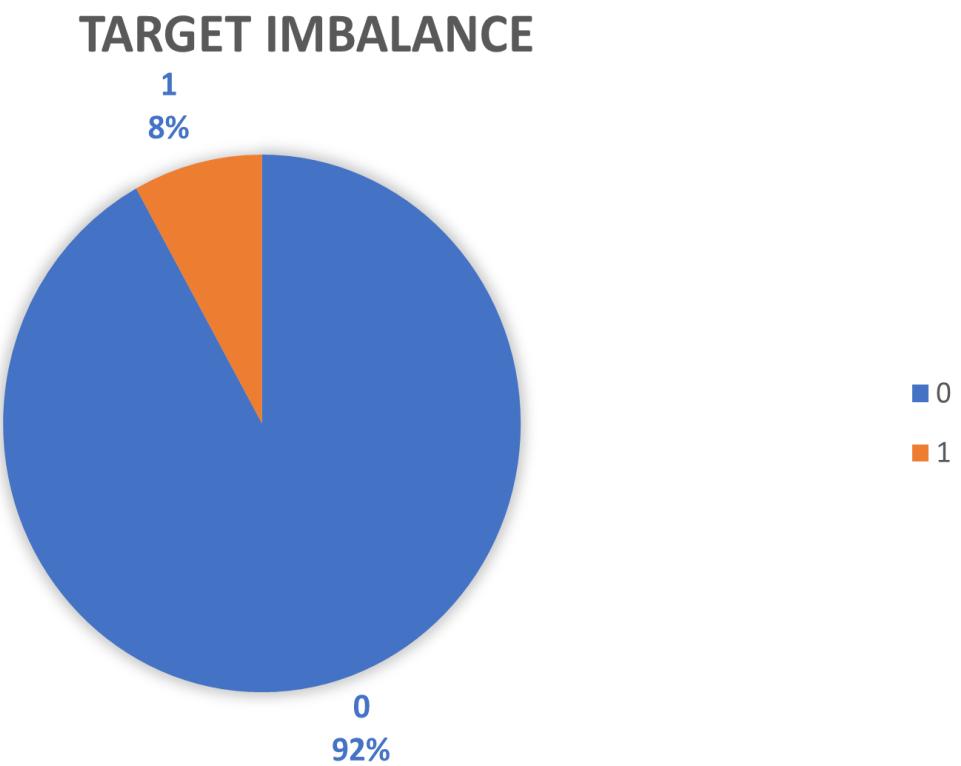
DAYS EMPLOYED



C. Analyze Data Imbalance:

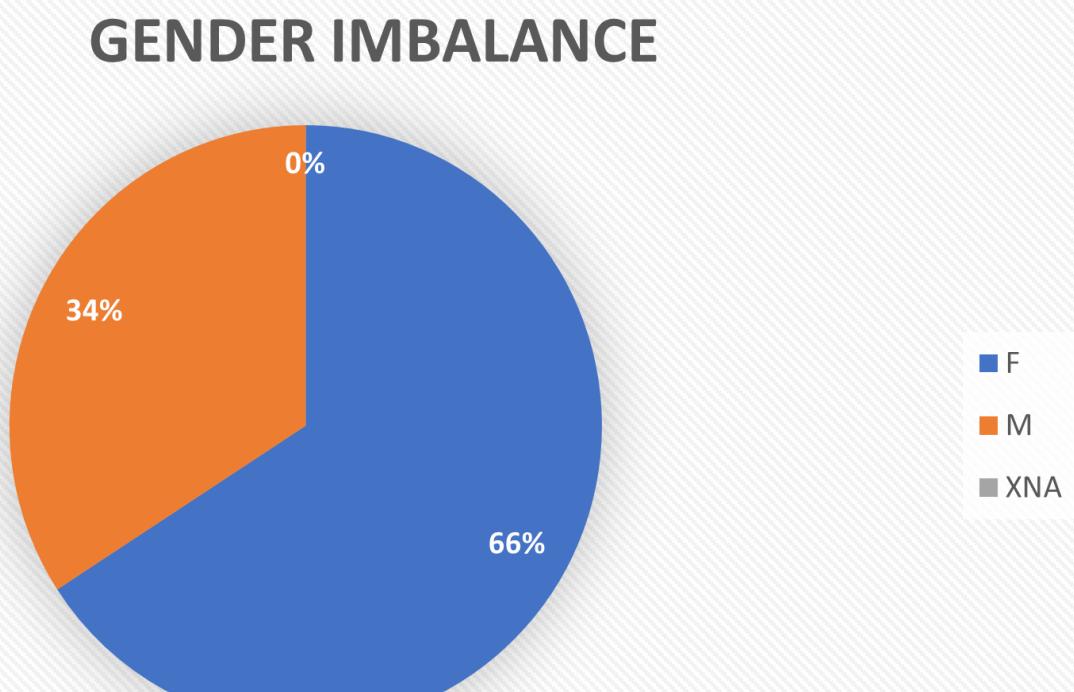
| Target Imbalance | Count of TARGET |
|--------------------|-----------------|
| 0 | 282686 |
| 1 | 24825 |
| Grand Total | 307511 |

The Target Imbalance Pie chart shows that almost 92% of the total clients had no problem during payment while 8% of the clients had some or the other problem.



| GENDER IMBALANCE | Count of CODE_GENDER |
|--------------------|----------------------|
| F | 202448 |
| M | 105059 |
| XNA | 4 |
| Grand Total | 307511 |

From the GENDER IMBALANCE pie chart we can infer that almost 66% of the clients are female and 34% of the clients are Male and the 4 of the applicants have gender as XNA which can be ignored.



| NAME OF THE HOUSING TYPE | Count of NAME OF THE HOUSING TYPE |
|----------------------------|-----------------------------------|
| Co-op apartment | 1122 |
| House / apartment | 272868 |
| Municipal apartment | 11183 |
| Office apartment | 2617 |
| Rented apartment | 4881 |
| With parents | 14840 |
| Grand Total | 307511 |

NAME OF THE HOUSING TYPE



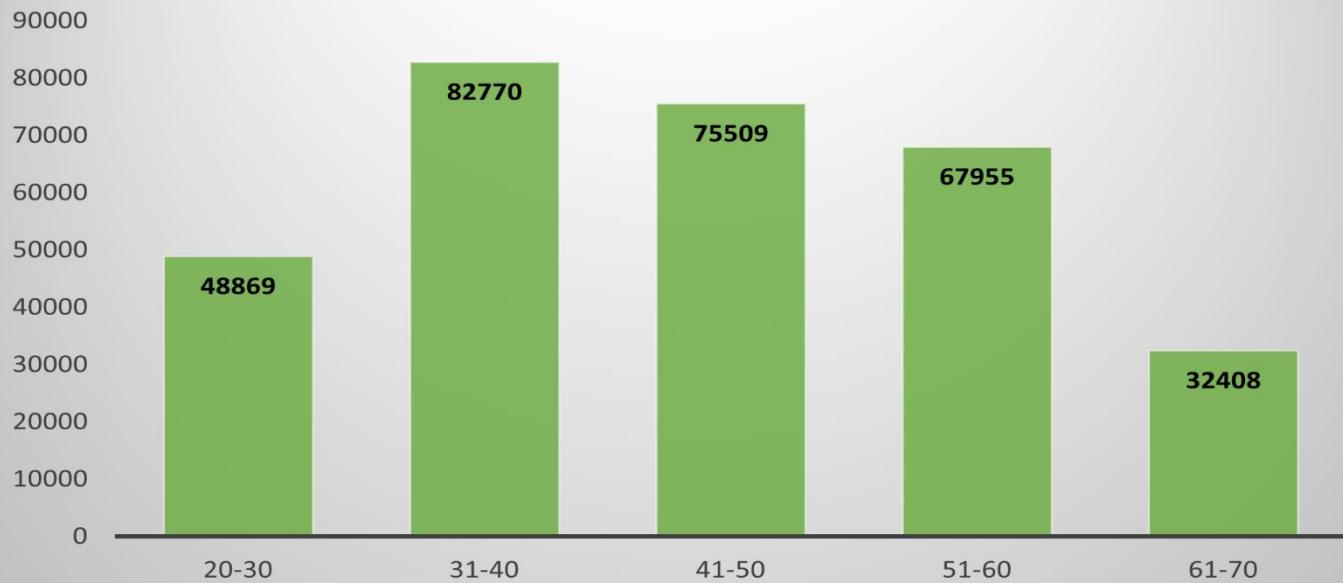
From the bar graph. The bank can target those groups who do not have their own apartment i.e. the bank may consider the people living in Co-op apartment, Municipal Apartment, Rented Apartment and people living with their parent.

D. Perform Univariate, Segmented Univariate, and Bivariate Analysis:

UNIVARIATE

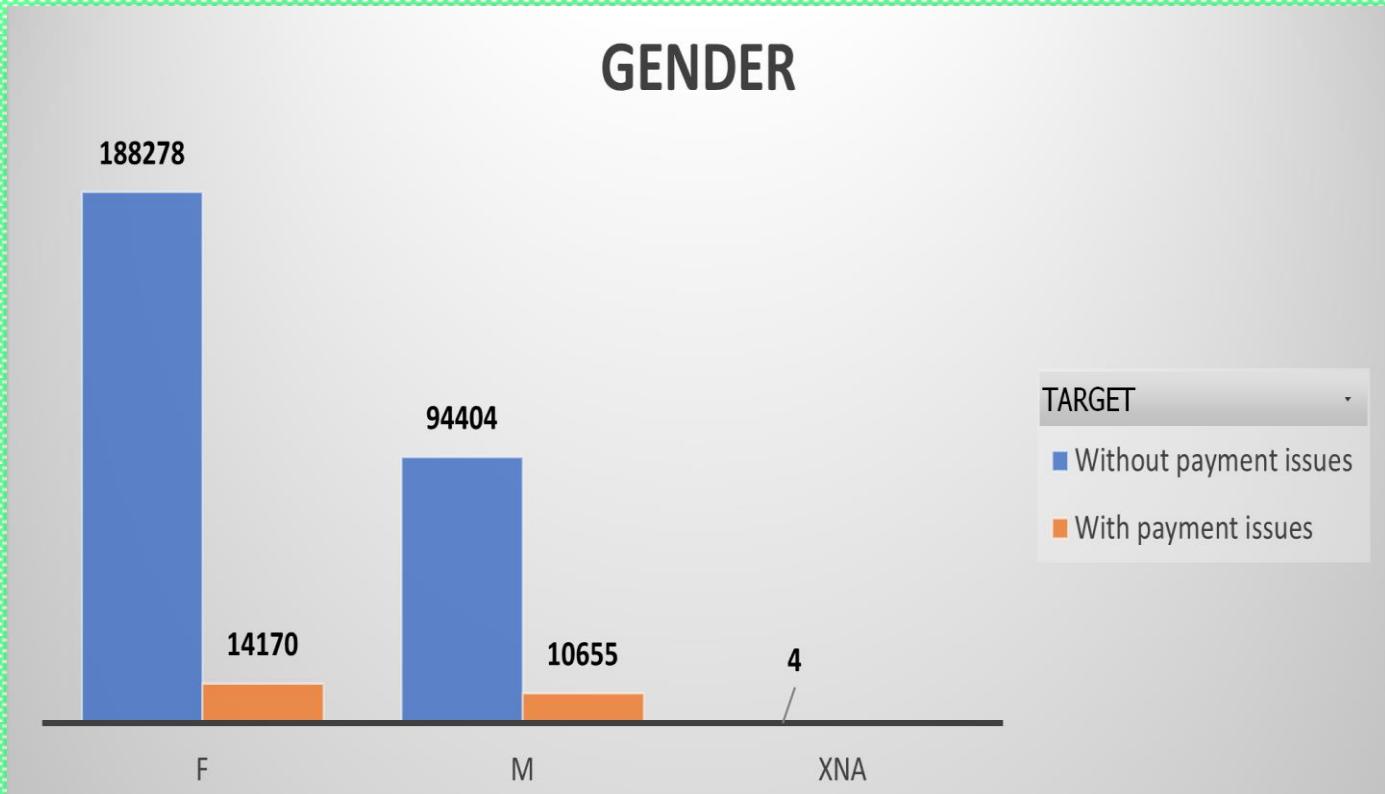
| AGE GROUPS | Count of YEARS BIRTH RANGE |
|--------------------|----------------------------|
| 20-30 | 48869 |
| 31-40 | 82770 |
| 41-50 | 75509 |
| 51-60 | 67955 |
| 61-70 | 32408 |
| Grand Total | 307511 |

AGE GROUP



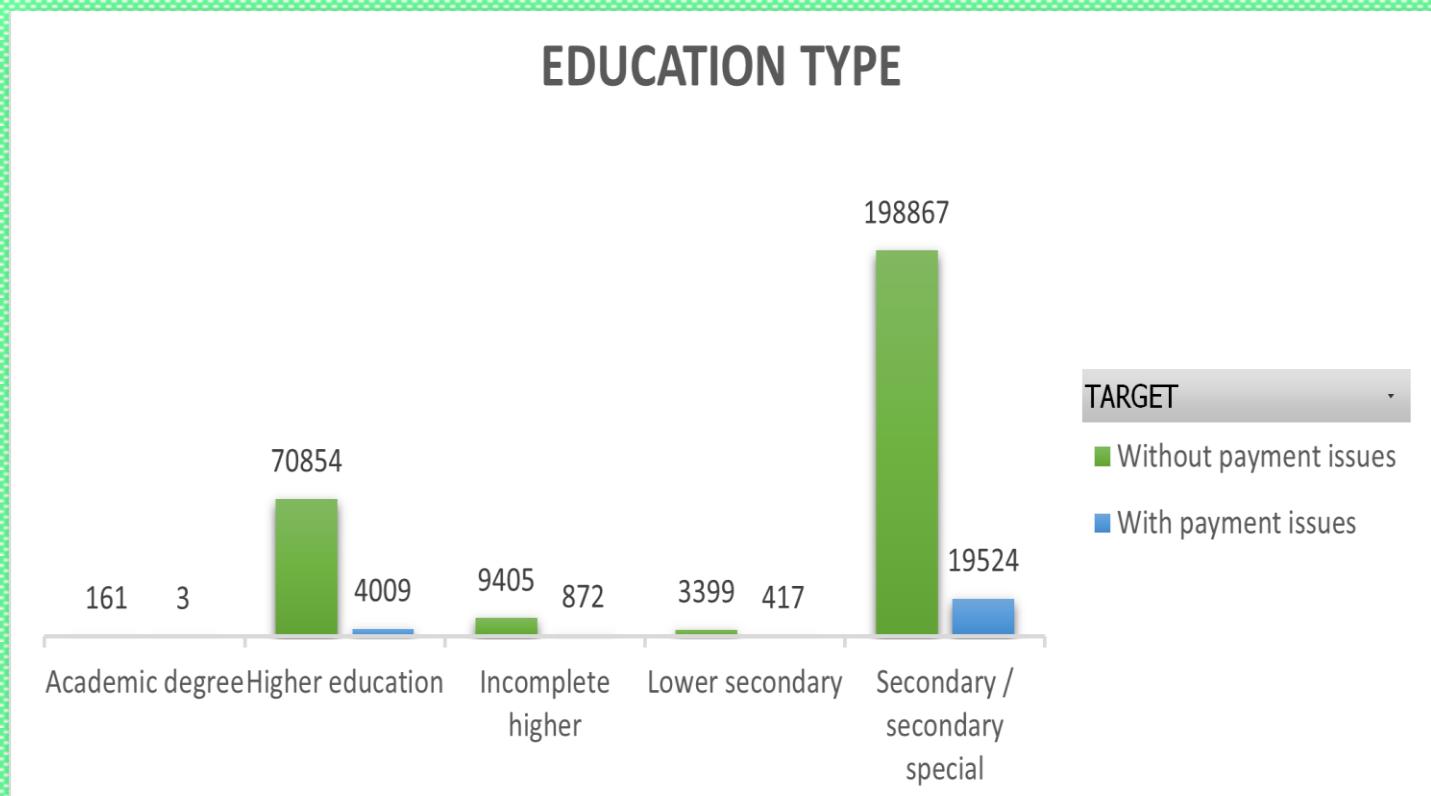
From the bar graph we can infer that most of the applicants belong to the **Age Group 31-40**.

| GENDER | Without payment issues | With payment issues | Grand Total |
|--------------------|------------------------|---------------------|---------------|
| F | 188278 | 14170 | 202448 |
| M | 94404 | 10655 | 105059 |
| XNA | 4 | / | 4 |
| Grand Total | 282686 | 24825 | 307511 |



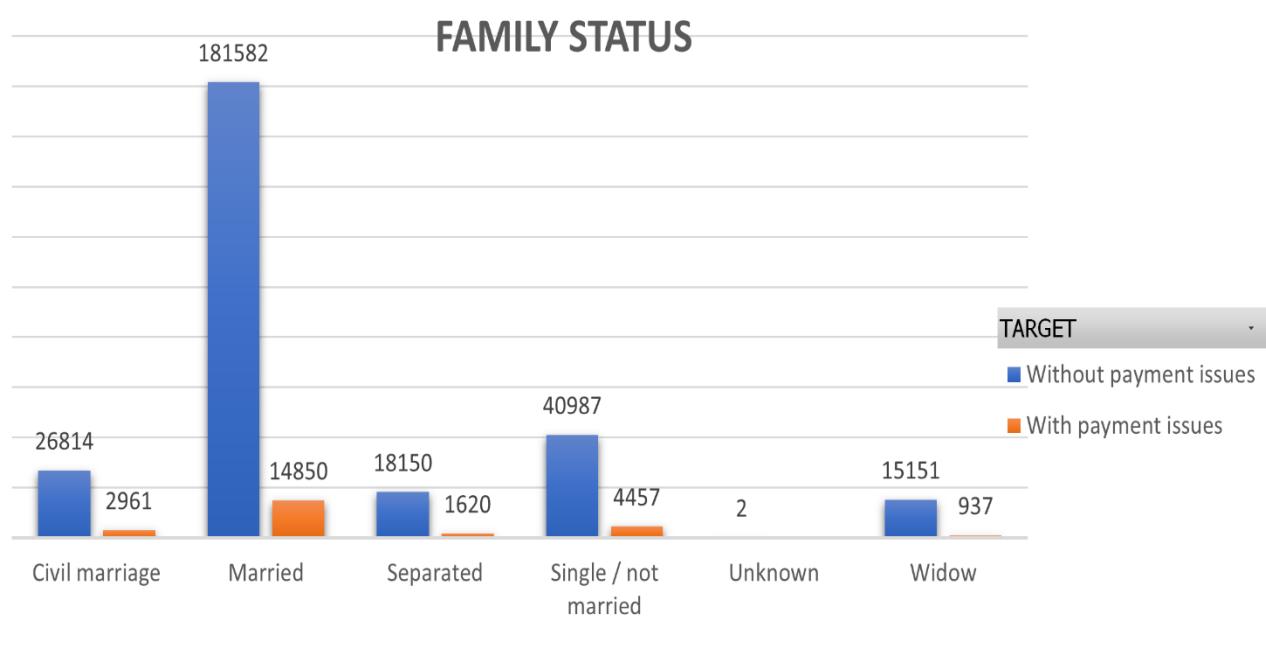
From the Bar graph we can infer that Clients with GENDER = 'F' have the highest number of non-defaulters i.e. $188278 - 14170 = \mathbf{174108}$.

| EDUCATION TYPE | Without payment issues | With payment issues | Grand Total |
|-------------------------------|------------------------|---------------------|---------------|
| Academic degree | 161 | 3 | 164 |
| Higher education | 70854 | 4009 | 74863 |
| Incomplete higher | 9405 | 872 | 10277 |
| Lower secondary | 3399 | 417 | 3816 |
| Secondary / secondary special | 198867 | 19524 | 218391 |
| Grand Total | 282686 | 24825 | 307511 |



From the above Bar graph we can infer that clients having EDUCATION TYPE = ‘SECONDARY/SECONDARY SPECIAL’ have the highest count for Non defaulters.

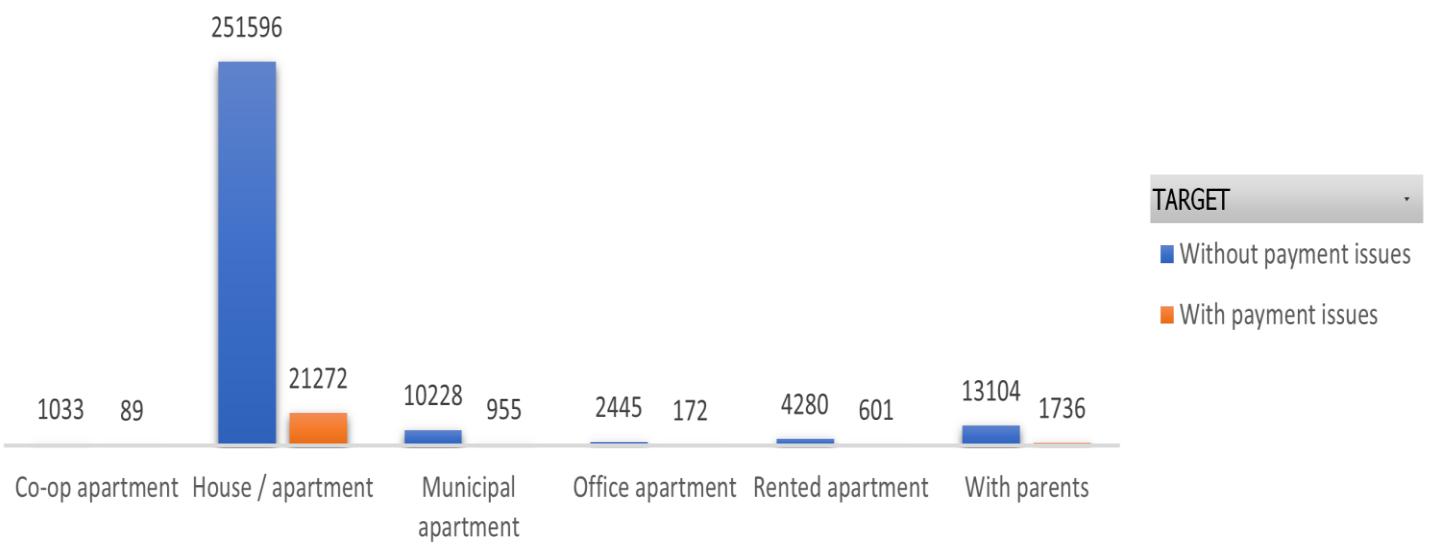
| FAMILY STATUS | Without payment issues | With payment issues | Grand Total |
|-----------------------------|-------------------------------|----------------------------|--------------------|
| Civil marriage | 26814 | 2961 | 29775 |
| Married | 181582 | 14850 | 196432 |
| Separated | 18150 | 1620 | 19770 |
| Single / not married | 40987 | 4457 | 45444 |
| Unknown | 2 | | 2 |
| Widow | 15151 | 937 | 16088 |
| Grand Total | 282686 | 24825 | 307511 |



From the above Bar chart we can infer that clients having FAMILY STATUS = ‘MARRIED’ have the highest count of Non defaulters.

| HOUSING TYPE | Without payment issues | With payment issues | Grand Total |
|----------------------------|-------------------------------|----------------------------|--------------------|
| Co-op apartment | 1033 | 89 | 1122 |
| House / apartment | 251596 | 21272 | 272868 |
| Municipal apartment | 10228 | 955 | 11183 |
| Office apartment | 2445 | 172 | 2617 |
| Rented apartment | 4280 | 601 | 4881 |
| With parents | 13104 | 1736 | 14840 |
| Grand Total | 282686 | 24825 | 307511 |

HOUSING TYPE

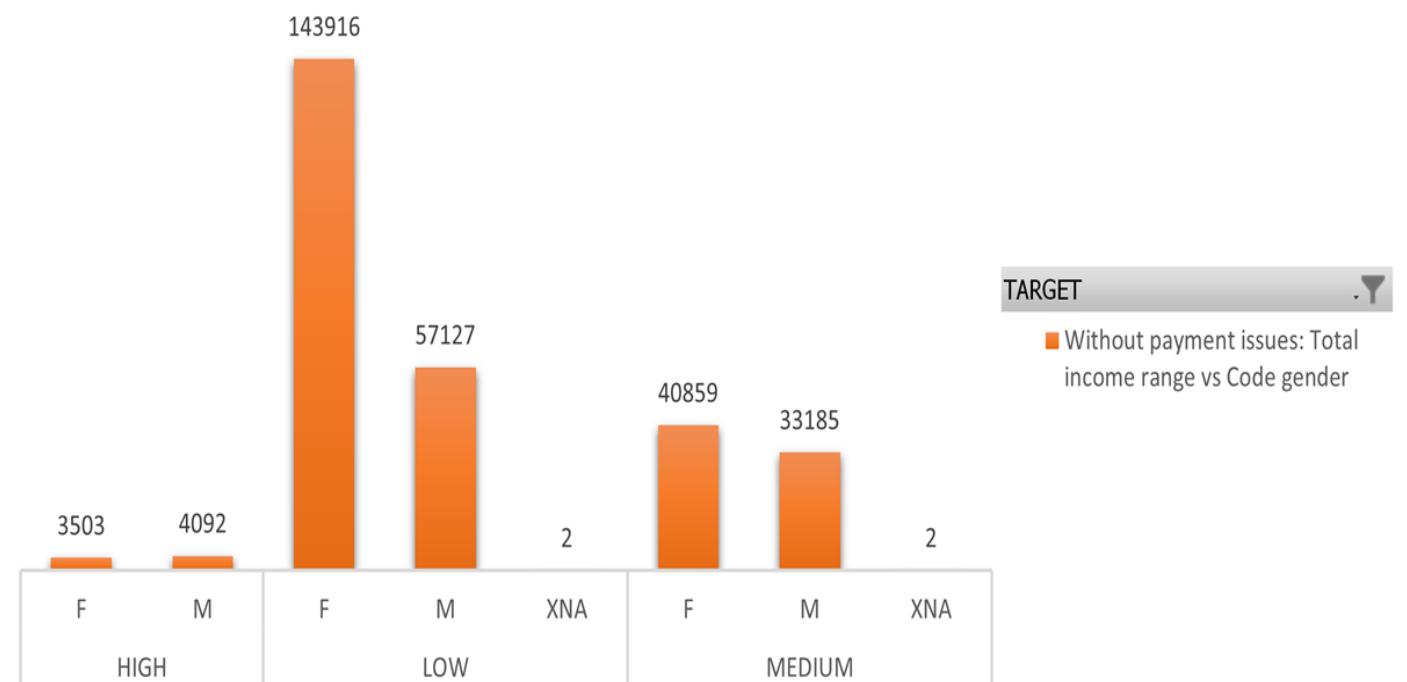


From the above Bar Chart we can infer that clients having **HOUSING TYPE = ‘House/Apartment’** have the highest count of Non-defaulters.

Bivariate

| Total income range and Code gender | Without payment issues: Total income range vs Code gender | Grand Total |
|------------------------------------|-----------------------------------------------------------|-------------|
| HIGH | 7595.00 | 7595.00 |
| F | 3503.00 | 3503.00 |
| M | 4092.00 | 4092.00 |
| LOW | 201045.00 | 201045.00 |
| F | 143916.00 | 143916.00 |
| M | 57127.00 | 57127.00 |
| XNA | 2.00 | 2.00 |
| MEDIUM | 74046.00 | 74046.00 |
| F | 40859.00 | 40859.00 |
| M | 33185.00 | 33185.00 |
| XNA | 2.00 | 2.00 |
| Grand Total | 282686.00 | 282686.00 |

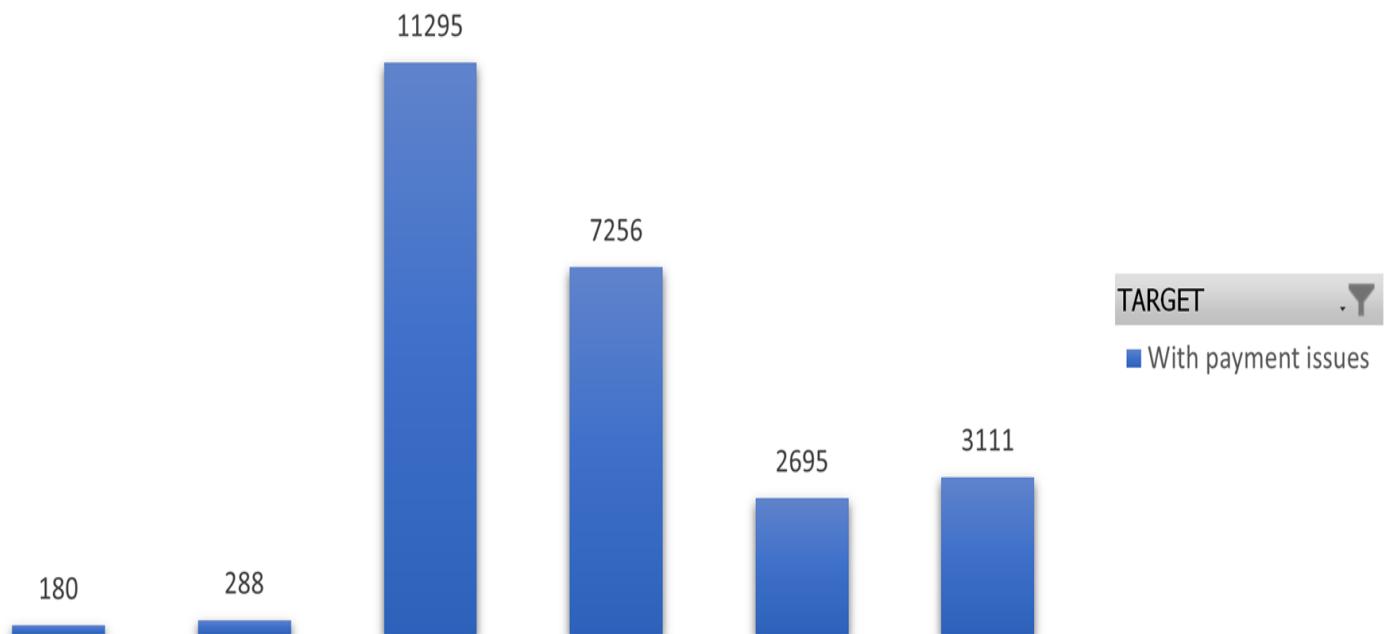
Without payment issues: Total income range vs Code gender



From the bar graph we can infer that Females belonging to Low income group are the highest number of clients with no payment issues.

| Total income range and Code gender | With payment issues | Grand Total |
|------------------------------------|---------------------|-------------|
| HIGH | 468 | 468 |
| F | 180 | 180 |
| M | 288 | 288 |
| LOW | 18551 | 18551 |
| F | 11295 | 11295 |
| M | 7256 | 7256 |
| MEDIUM | 5806 | 5806 |
| F | 2695 | 2695 |
| M | 3111 | 3111 |
| Grand Total | 24825 | 24825 |

Total income range vs Code gender



From the bar graph we can infer that Females belonging to Low income group are the highest number of clients with payment issues.

Previous Application Dataset:

The following columns of the previous application datasets need to be dropped as they are irrelevant for the analysis:

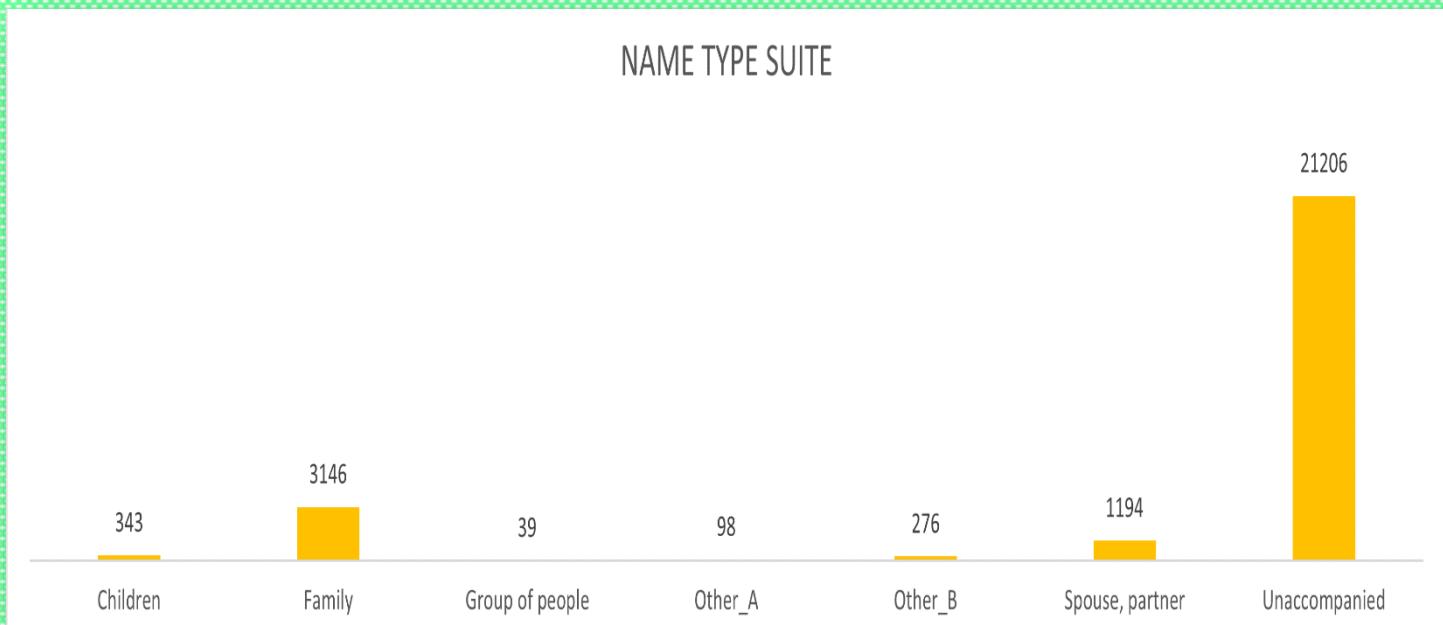
- HOUR_APPR_PROCESS_START
- WEEKDAY_APPR_PROCESS_START_PREV
- FLAG_LAST_APPL_PER_CONTRACT
- NFLAG_LAST_APPL_IN_DAY
- SK_ID_CURR
- WEEKDAY_APPR_PROCESS_START

Removing the rows with the values 'XNA' &'XAP' from the column:

NAME_TYPE_SUITE

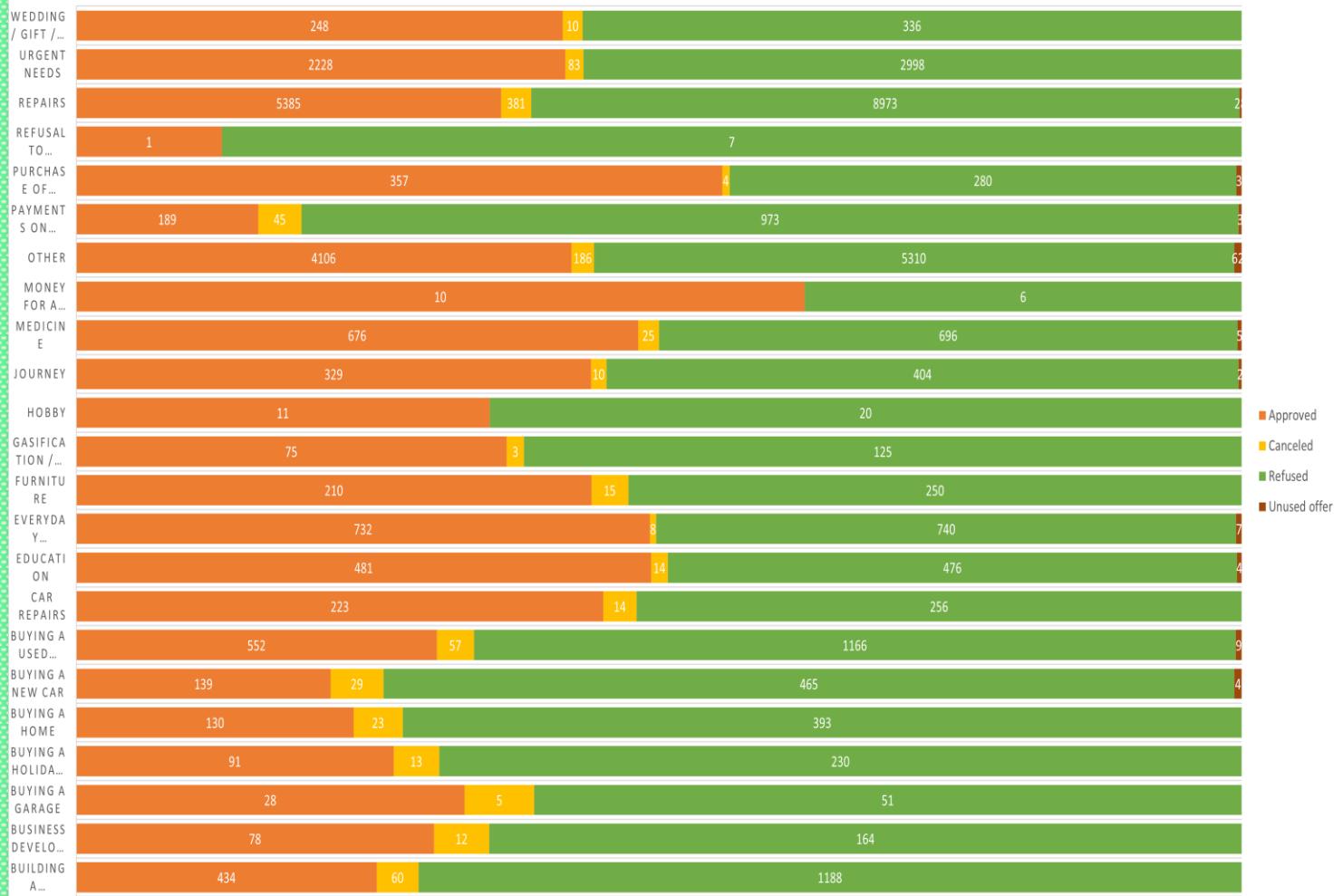
In the column **AMT_ANNUITY** replace the blank rows with **Median of AMT_ANNUITY '21340'**.

| NAME_TYPE_SUITE | Count of NAME_TYPE_SUITE |
|--------------------|--------------------------|
| Children | 343 |
| Family | 3146 |
| Group of people | 39 |
| Other_A | 98 |
| Other_B | 276 |
| Spouse, partner | 1194 |
| Unaccompanied | 21206 |
| Grand Total | 26302 |



From the bar graph we can infer that Unaccompanied has the highest count therefore replace blank rows with Unaccompanied.

DISTRIBUTION OF NAME CONTRACT STATUS



| Name Contract Status | Approved | Canceled | Refused | Unused offer | Grand Total |
|----------------------------------|--------------|----------|------------|--------------|--------------|
| Building a house or an annex | 434 | | 60 | 1188 | 1682 |
| Business development | 78 | | 12 | 164 | 254 |
| Buying a garage | 28 | | 5 | 51 | 84 |
| Buying a holiday home / land | 91 | | 13 | 230 | 334 |
| Buying a home | 130 | | 23 | 393 | 546 |
| Buying a new car | 139 | | 29 | 465 | 637 |
| Buying a used car | 552 | | 57 | 1166 | 1784 |
| Car repairs | 223 | | 14 | 256 | 493 |
| Education | 481 | | 14 | 476 | 975 |
| Everyday expenses | 732 | | 8 | 740 | 1487 |
| Furniture | 210 | | 15 | 250 | 475 |
| Gasification / water supply | 75 | | 3 | 125 | 203 |
| Hobby | 11 | | | 20 | 31 |
| Journey | 329 | | 10 | 404 | 745 |
| Medicine | 676 | | 25 | 696 | 1402 |
| Money for a third person | 10 | | | 6 | 16 |
| Other | 4106 | | 186 | 5310 | 9664 |
| Payments on other loans | 189 | | 45 | 973 | 1210 |
| Purchase of electronic equipment | 357 | | 4 | 280 | 644 |
| Refusal to name the goal | 1 | | | 7 | 8 |
| Repairs | 5385 | | 381 | 8973 | 14767 |
| Urgent needs | 2228 | | 83 | 2998 | 5309 |
| Wedding / gift / holiday | 248 | | 10 | 336 | 594 |
| Grand Total | 16713 | | 997 | 25507 | 43344 |

From the graph and table we can infer that Name of Contract status i.e. **Repairs** work has the highest count of Approved and Refused Loans.

Conclusion

From the above analysis, we can infer that what kind of peoples can repay loan, what kinds of loan people prefer to take, people taking loans come from which background, what is their source of income, for what type of people the loan applications are refused and based on which conditions.

Impact of Car Features on

Price and Profitability

Project-7



Project Description:

The aim of this project is to analyze the impact of car features on price and profitability in the automotive industry. By analyzing the dataset containing information on various car models and their specifications and their prices.

Approach:

Firstly I have cleaned the dataset to ensure accurate and reliable results for the analysis. Then I have performed analytical methods and Visualization Methods. These techniques will help to get insights about the Impact of Car Features on Price and Profitability.

Tech-Stack Used:

PPT – To prepare a detailed report.

EXCEL – Excel was used to perform entire analysis.

Insight:

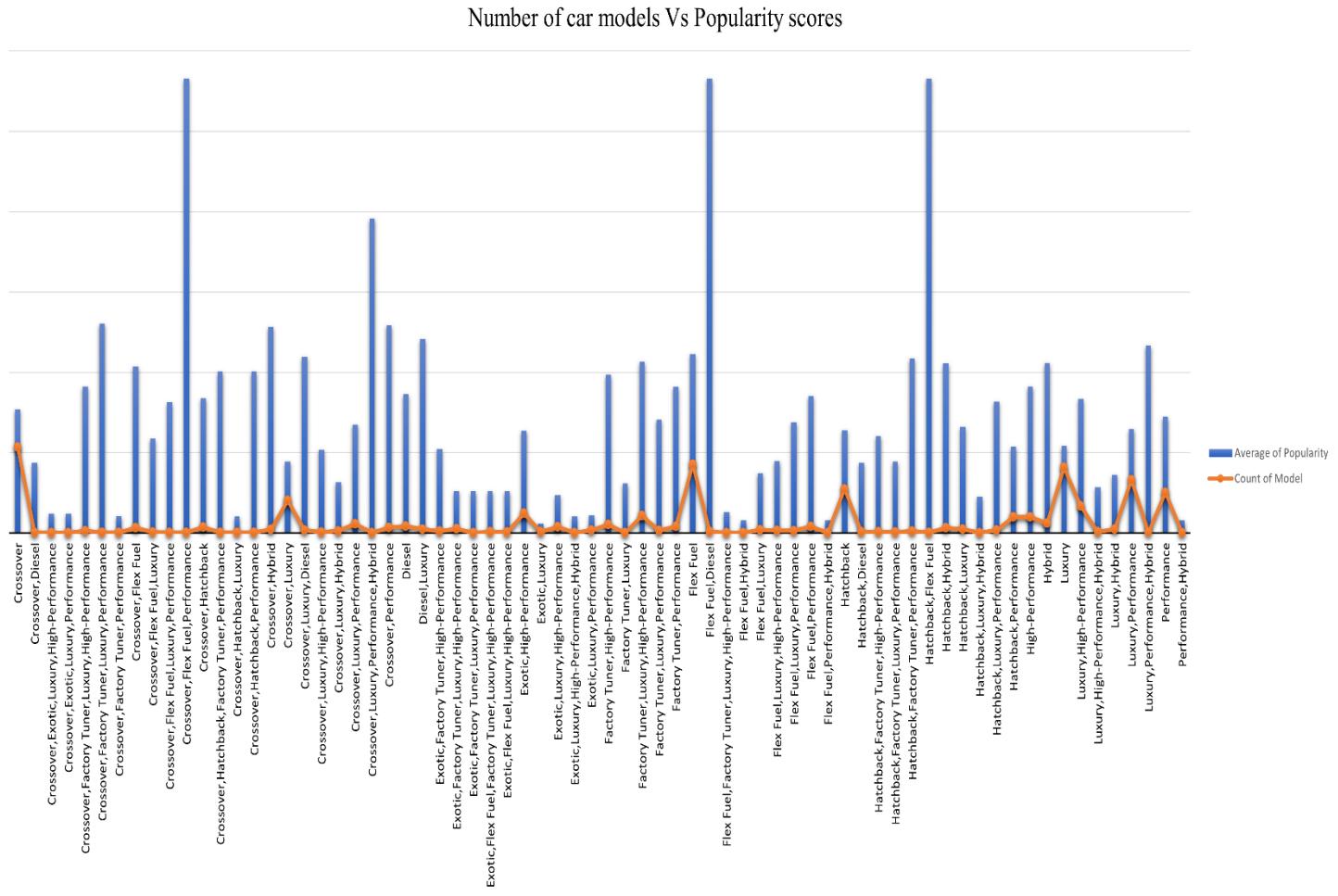
1.How does the popularity of a car model vary across different market categories?

Task 1.A: Create a pivot table that shows the number of car models in each market category and their corresponding popularity scores.

Pivot Table have been created to show the number of car models in each market category and their corresponding popularity scores.

Task 1.B: Create a combo chart that visualizes the relationship between market category and popularity.

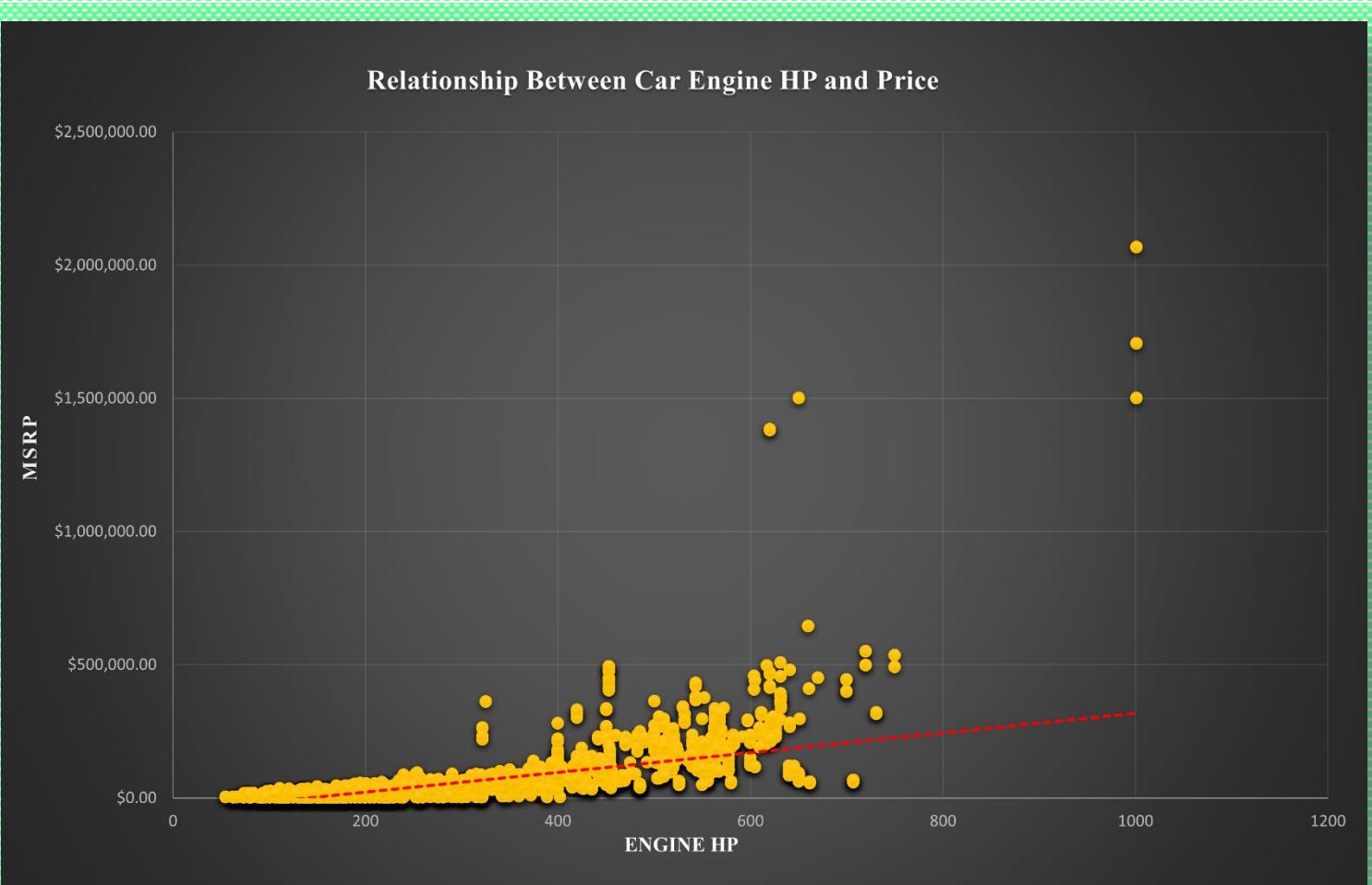
Chart have been created to show the relationship between market category and popularity. From the chart below we can infer that crossover, flex fuel, diesel, hatchback, and performance are the most popular market categories for car models.



2.What is the relationship between a car's engine power and its price?

Task 2: Create a scatter chart that plots engine power on the x-axis and price on the y-axis. Add a trendline to the chart to visualize the relationship between these variables.

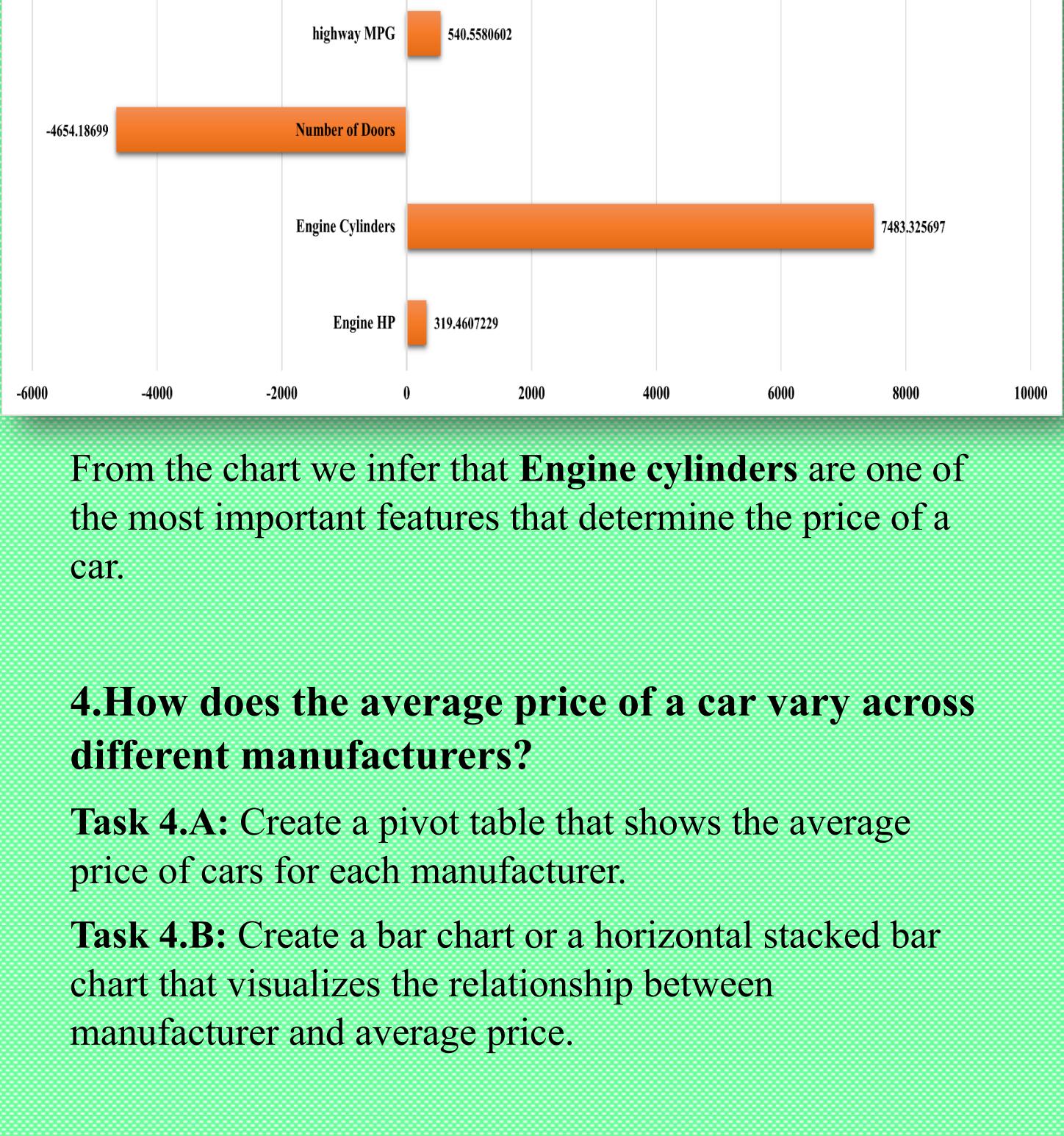
Scatter chart have been created in following slide from the chart we can infer that the horse power of the engine increases then the price also increases.



3.Which car features are most important in determining a car's price?

Task 3: Use regression analysis to identify the variables that have the strongest relationship with a car's price.
Then create a bar chart that shows the coefficient values

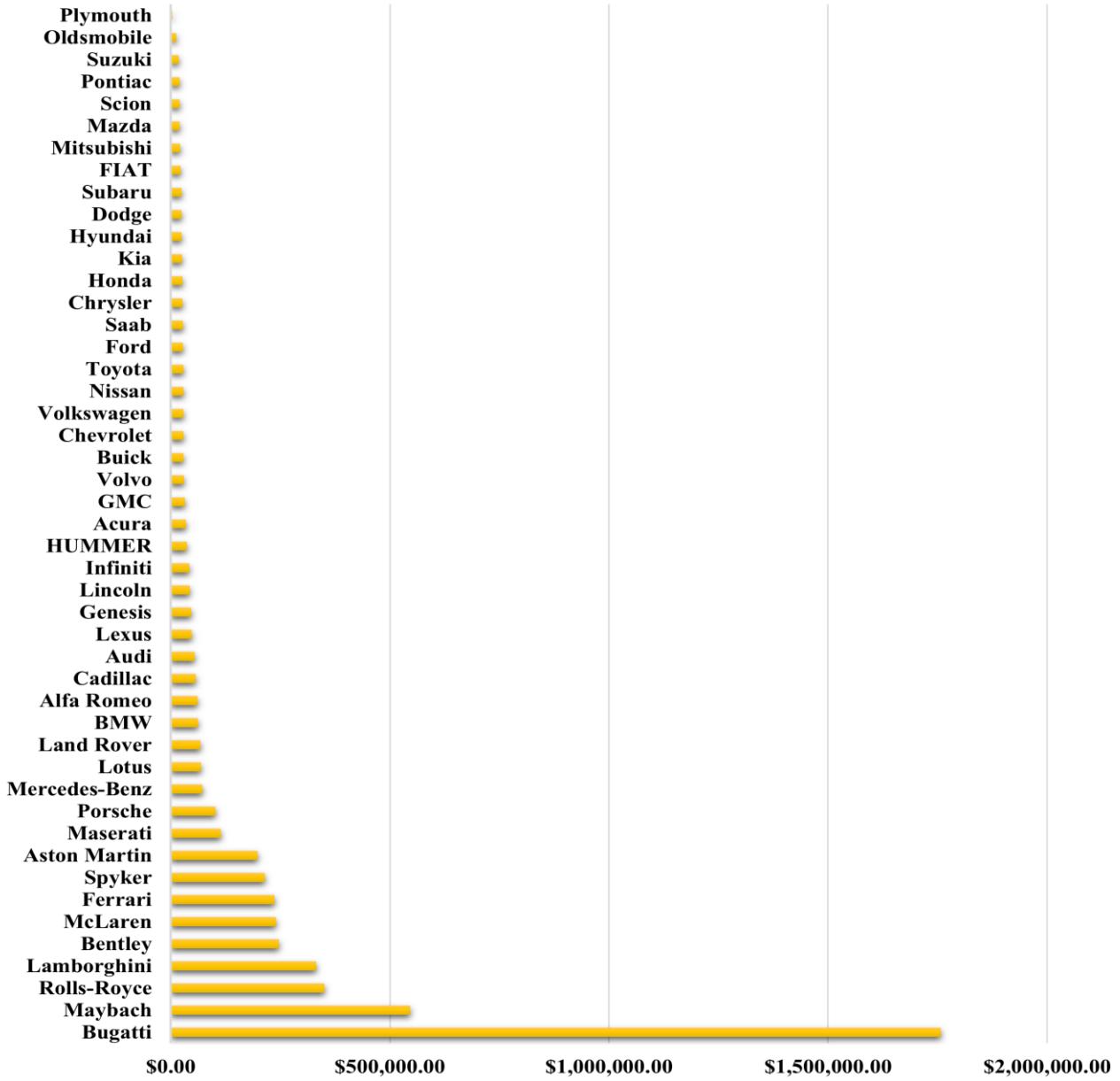
| Variables | Coefficients |
|-------------------------|--------------|
| Engine HP | 319.4607229 |
| Engine Cylinders | 7483.325697 |
| Number of Doors | -4654.18699 |
| highway MPG | 540.5580602 |
| city mpg | 1193.47926 |



| Manufacturers | Average Price |
|---------------|----------------|
| Bugatti | \$1,757,223.67 |
| Maybach | \$546,221.88 |
| Rolls-Royce | \$351,130.65 |
| Lamborghini | \$331,567.31 |
| Bentley | \$247,169.32 |
| McLaren | \$239,805.00 |
| Ferrari | \$237,383.82 |
| Spyker | \$214,990.00 |
| Aston Martin | \$198,123.46 |
| Maserati | \$113,684.49 |
| Porsche | \$101,622.40 |
| Mercedes-Benz | \$72,135.03 |
| Lotus | \$68,377.14 |
| Land Rover | \$68,067.09 |
| BMW | \$62,162.56 |
| Alfa Romeo | \$61,600.00 |
| Cadillac | \$56,368.27 |
| Audi | \$54,574.12 |
| Lexus | \$47,549.07 |
| Genesis | \$46,616.67 |
| Lincoln | \$43,560.01 |
| Infiniti | \$42,640.27 |

| Manufacturers | Average Price |
|---------------|---------------|
| HUMMER | \$36,464.41 |
| Acura | \$35,087.49 |
| GMC | \$32,444.09 |
| Volvo | \$29,724.68 |
| Buick | \$29,034.19 |
| Chevrolet | \$29,000.22 |
| Volkswagen | \$28,947.37 |
| Nissan | \$28,856.42 |
| Toyota | \$28,758.77 |
| Ford | \$28,522.86 |
| Saab | \$27,879.81 |
| Chrysler | \$26,722.96 |
| Honda | \$26,608.88 |
| Kia | \$25,318.75 |
| Hyundai | \$24,926.26 |
| Dodge | \$24,857.05 |
| Subaru | \$24,240.67 |
| FIAT | \$22,206.02 |
| Mitsubishi | \$21,316.35 |
| Mazda | \$20,106.56 |
| Scion | \$19,932.50 |
| Pontiac | \$19,800.04 |
| Suzuki | \$18,021.05 |
| Oldsmobile | \$12,843.80 |
| Plymouth | \$3,296.87 |

Manufacturer Vs Average price



From the chart we can infer that **Bugatti** has the highest average price.

5.What is the relationship between fuel efficiency and the number of cylinders in a car's engine?

Task 5. A: Create a scatter plot with the number of cylinders on the x-axis and highway MPG on the y-axis. Then create a trendline on the scatter plot to visually estimate the slope of the relationship and assess its significance.

Task 5. B: Calculate the correlation coefficient between the number of cylinders and highway MPG to quantify the strength and direction of the relationship.

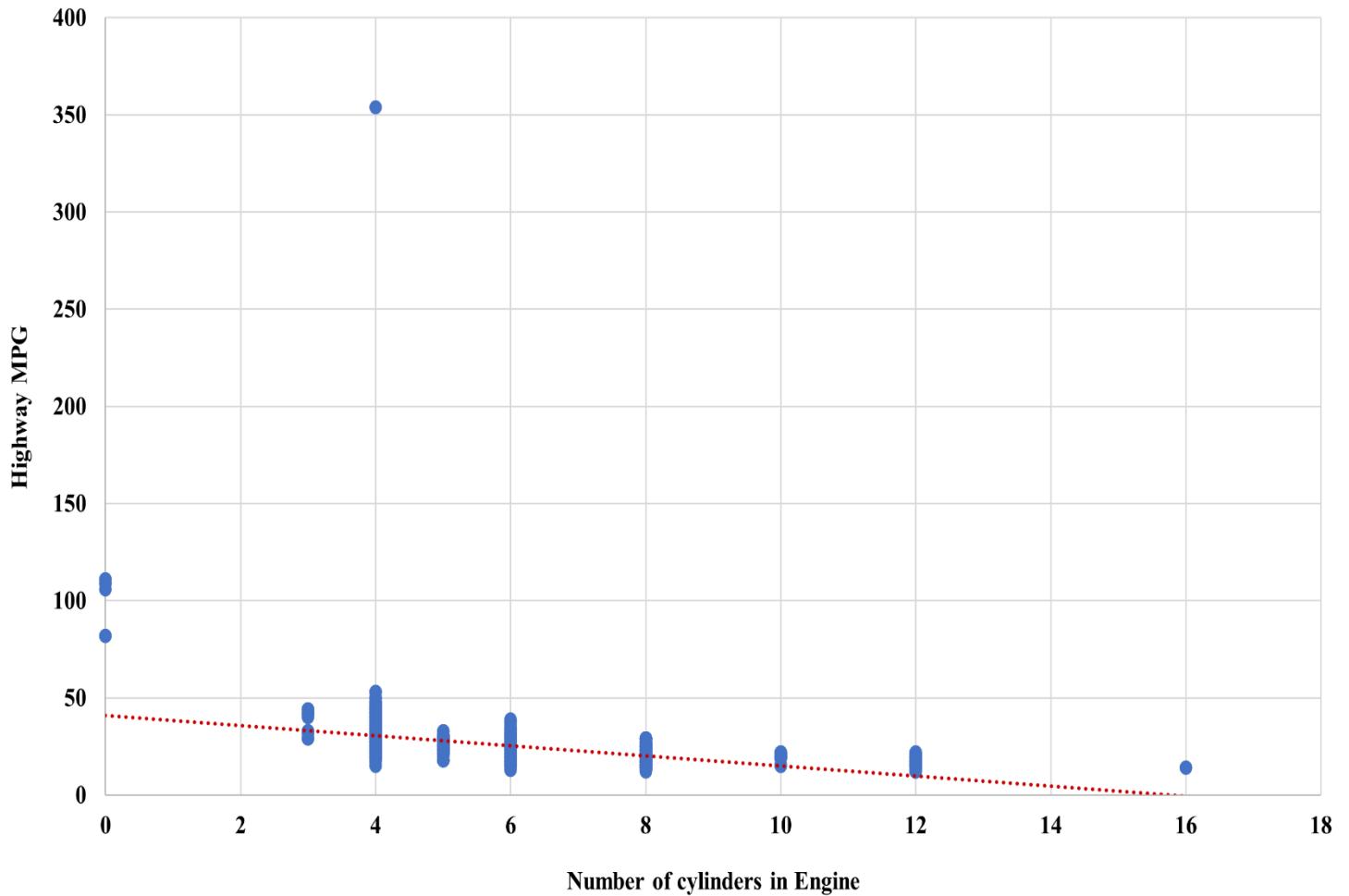
A. Created a **scatter chart** with the number of cylinders on the x-axis and highway MPG on the y-axis and a trendline have been plotted.

B. **Correlation coefficient** between the number of cylinders and highway MPG to quantify the strength and direction of the relationship have been calculated using the formula:

=CORREL(A2:A11098,B2:B11098)

| | |
|-------------------------|-----------|
| Correlation Coefficient | -0.614703 |
|-------------------------|-----------|

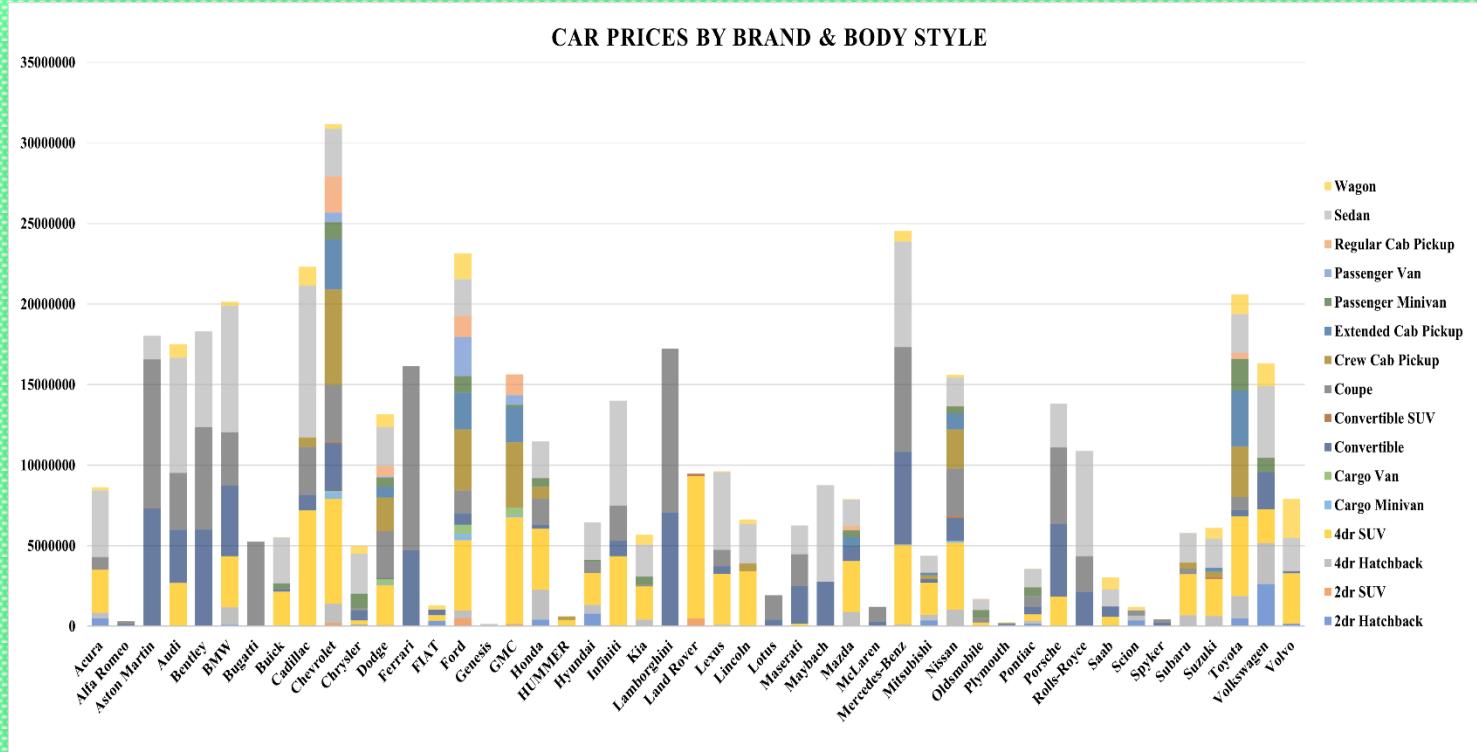
Number of cylinders Vs Highway MPG



From the chart we can infer that if the **number of cylinders increases then the highway mpg will decrease.**

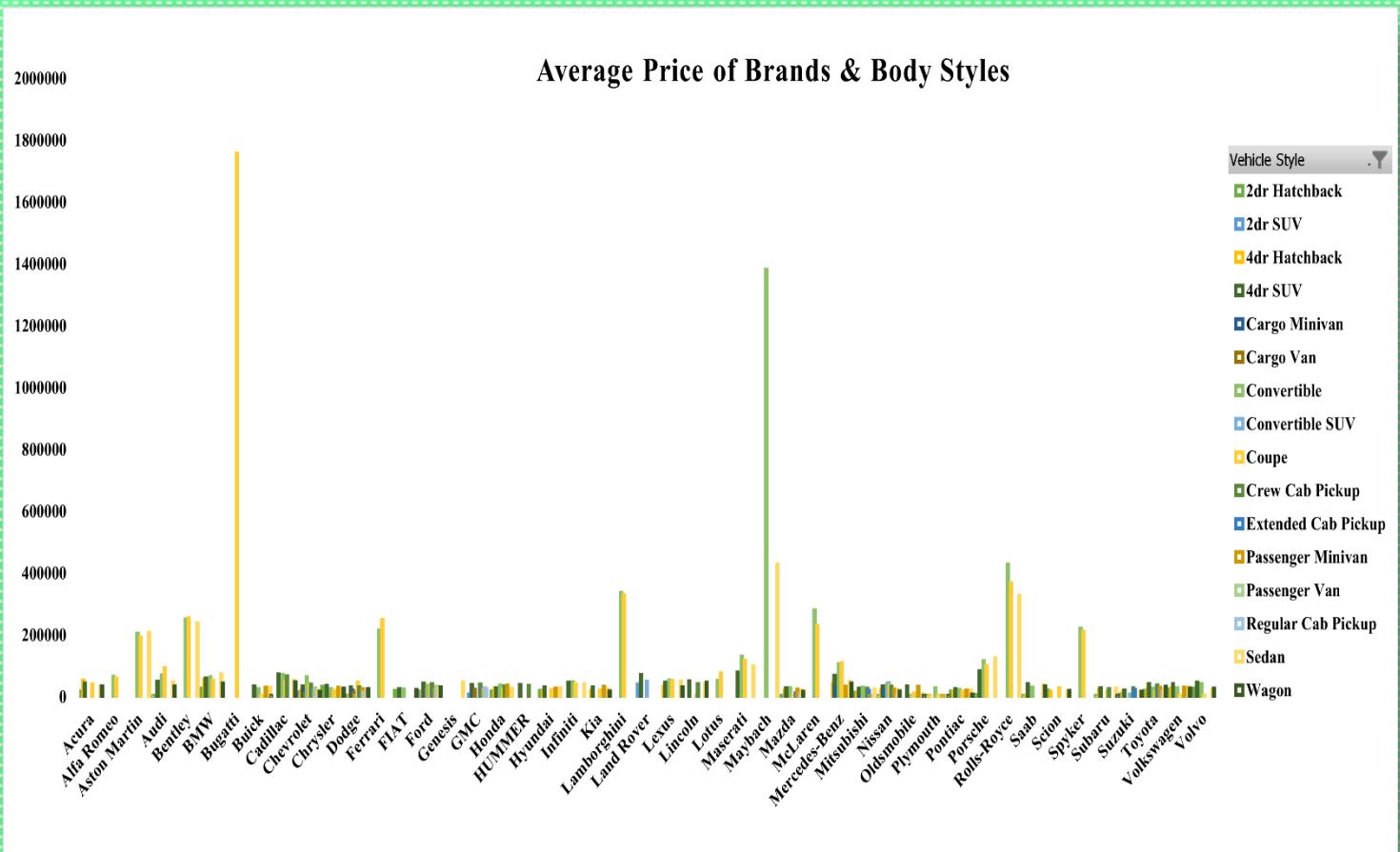
Building the Dashboard:

Task 1: How does the distribution of car prices vary by brand and body style?



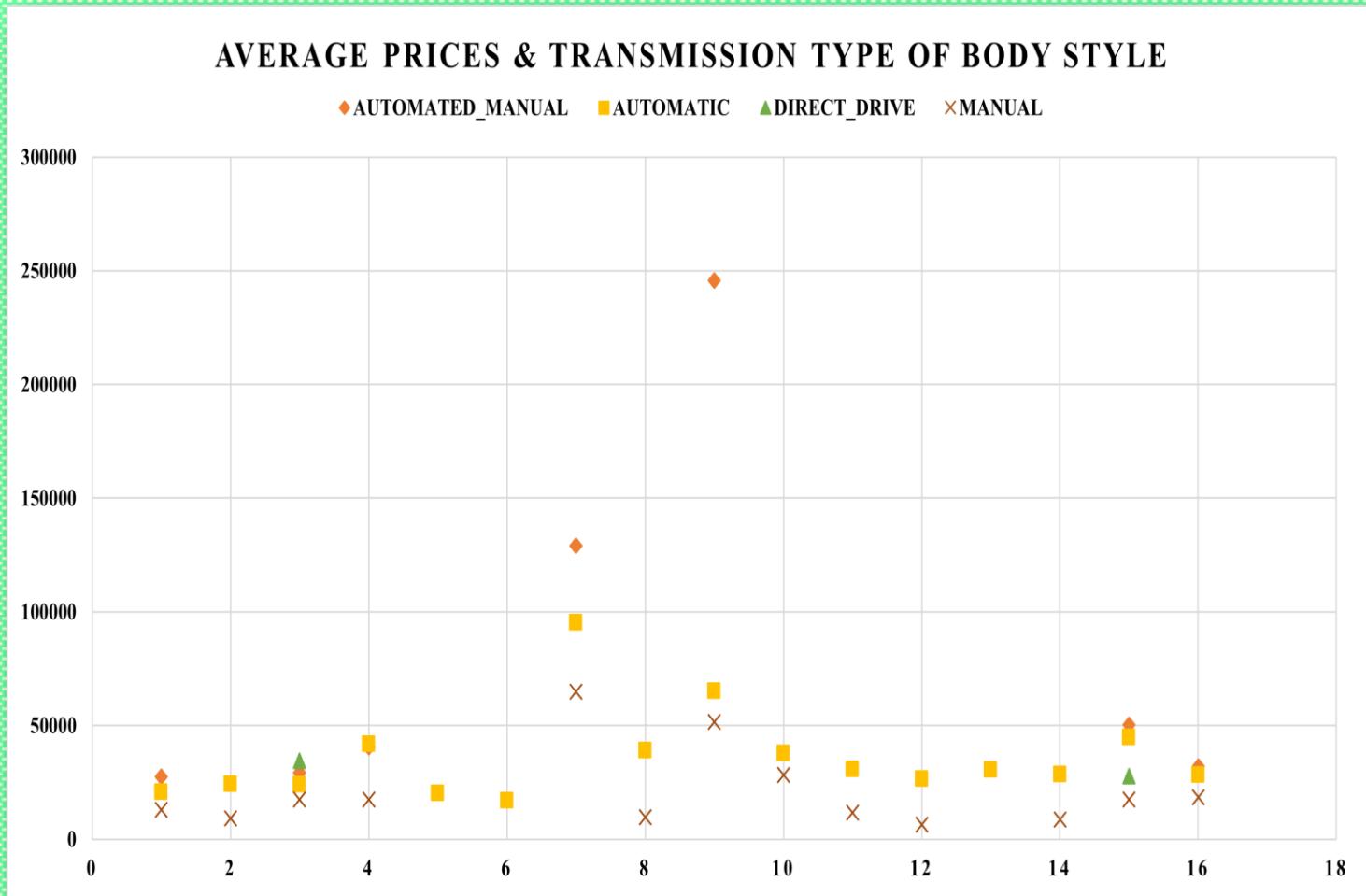
From the chart we can infer that **Chevrolet** has the highest price distribution by body style.

Task 2: Which car brands have the highest and lowest average MSRPs, and how does this vary by body style?



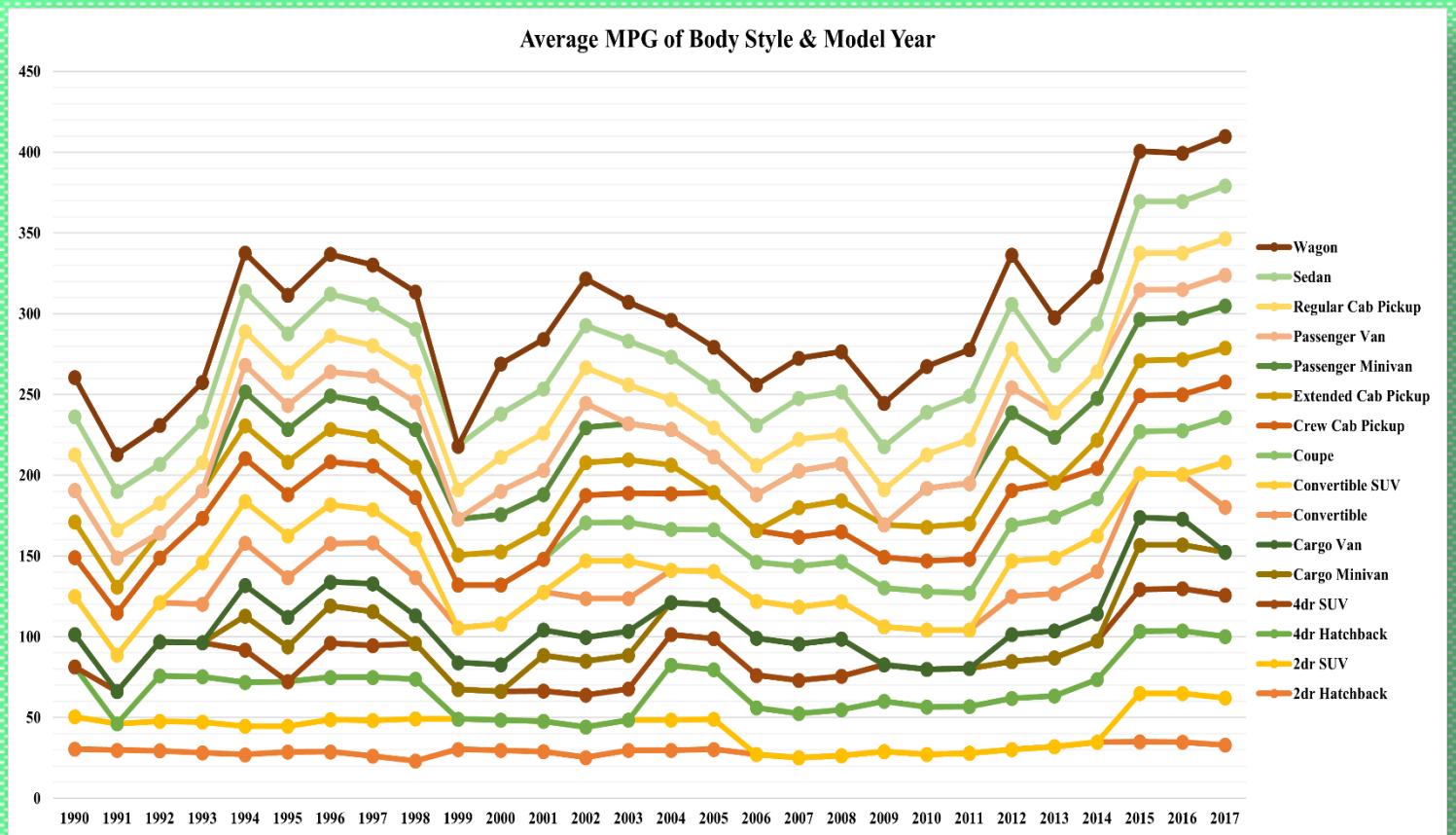
From the chart we can infer that **Bugatti has the highest average price** and **Plymouth has the lowest average price.**

Task 3: How do the different feature such as transmission type affect the MSRP, and how does this vary by body style?



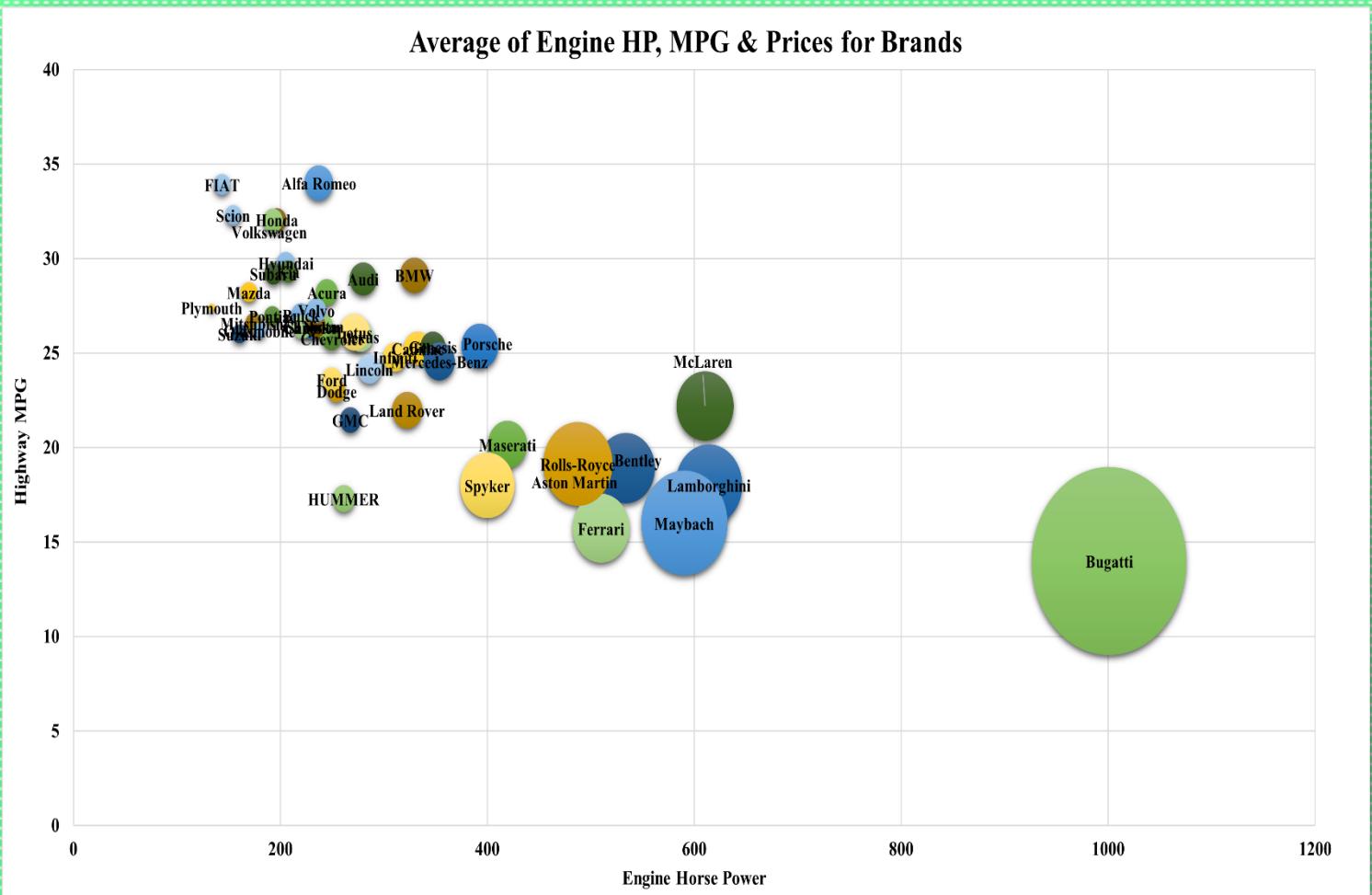
From the chart we can infer that **automated manual** is the most expensive transmission.

Task 4: How does the fuel efficiency of cars vary across different body styles and model years?



From the chart we can infer that fuel efficiency of cars increased across different body styles.

Task 5: How does the car's horsepower, MPG, and price vary across different Brands?



From the chart we can infer that if the engine horse power increases then highway mpg will decrease and the price also will increase hence making Bugatti more expensive.

Conclusion

From the above analysis, we can infer that the provided dataset has several meaningful insights which can help the car manufacturers to optimize pricing and product development decisions to maximize profitability while meeting consumer demand.

ABC CALL VOLUME

TREND ANALYSIS

PROJECT-8



Project Description:

In this project I have provided with a dataset of a Customer Experience (CX) analytics, specifically focusing on the inbound calling team of a company. You'll be provided with a dataset that spans 23 days and includes various details and the goal is to attract, engage, and delight customers, turning them into loyal advocates for the business. Additionally, the project involves addressing the issue of unanswered calls during the night and proposing a manpower plan for that time period as well.

Approach:

Firstly I have cleaned the dataset to ensure accurate and reliable results for the analysis. Then I have performed analytical methods and Visualization Methods. These techniques will help to get insights about the ABC Call Volume Trend Analysis.

Tech-Stack Used:

PPT – To prepare a detailed report.

EXCEL – Excel was used to perform entire analysis.

Insights:

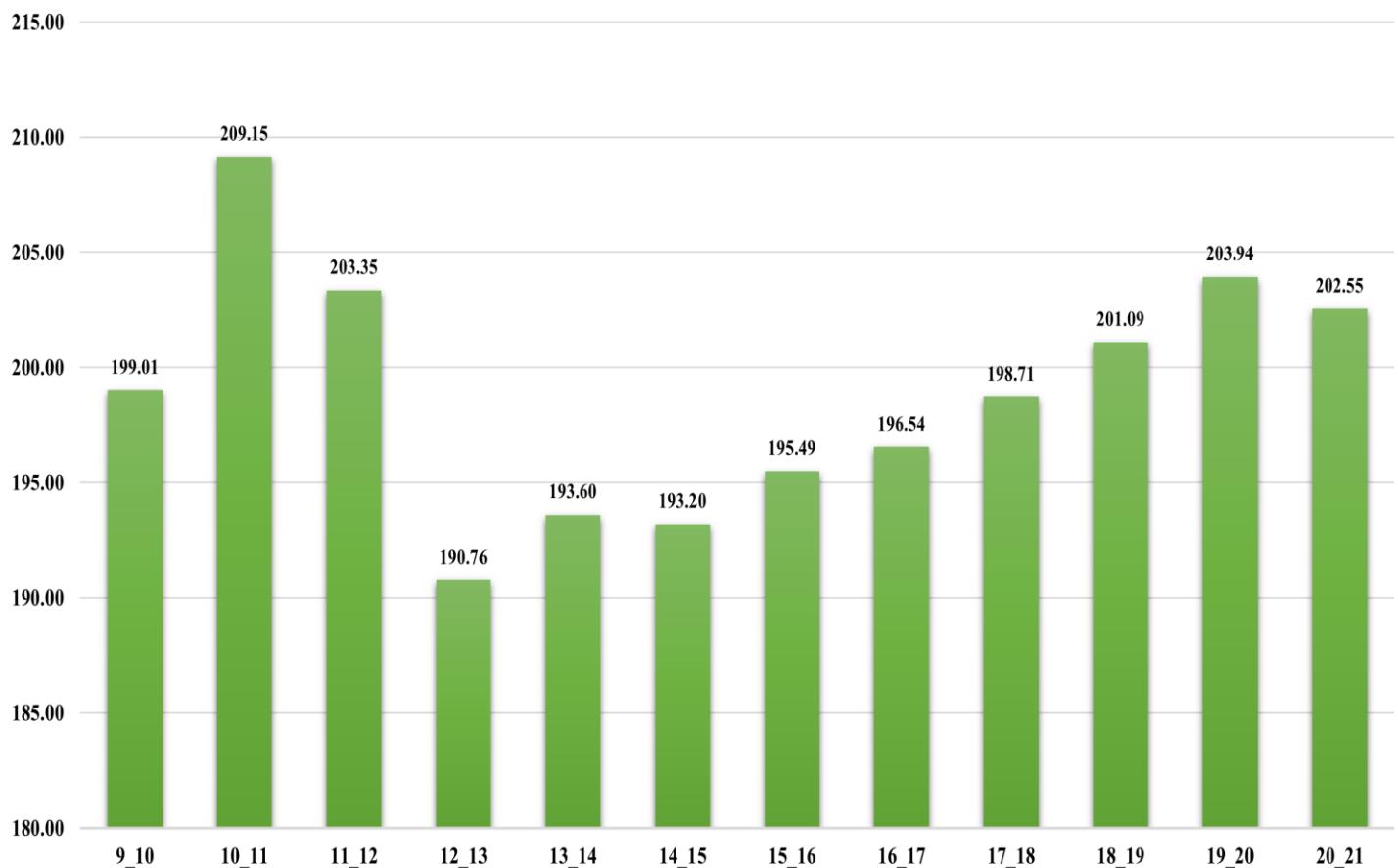
1. Average Call Duration:

What is the average duration of calls for each time bucket?

From the above bar chart we can infer that time bucket 10_11 i.e. 10AM to 11AM has the highest of average of calls in seconds i.e. **209.15**

| Time bucket | Average of Call Seconds (s) |
|-------------|-----------------------------|
| 9_10 | 199.01 |
| 10_11 | 209.15 |
| 11_12 | 203.35 |
| 12_13 | 190.76 |
| 13_14 | 193.60 |
| 14_15 | 193.20 |
| 15_16 | 195.49 |
| 16_17 | 196.54 |
| 17_18 | 198.71 |
| 18_19 | 201.09 |
| 19_20 | 203.94 |
| 20_21 | 202.55 |

Average duration of all calls (in sec) received by agents

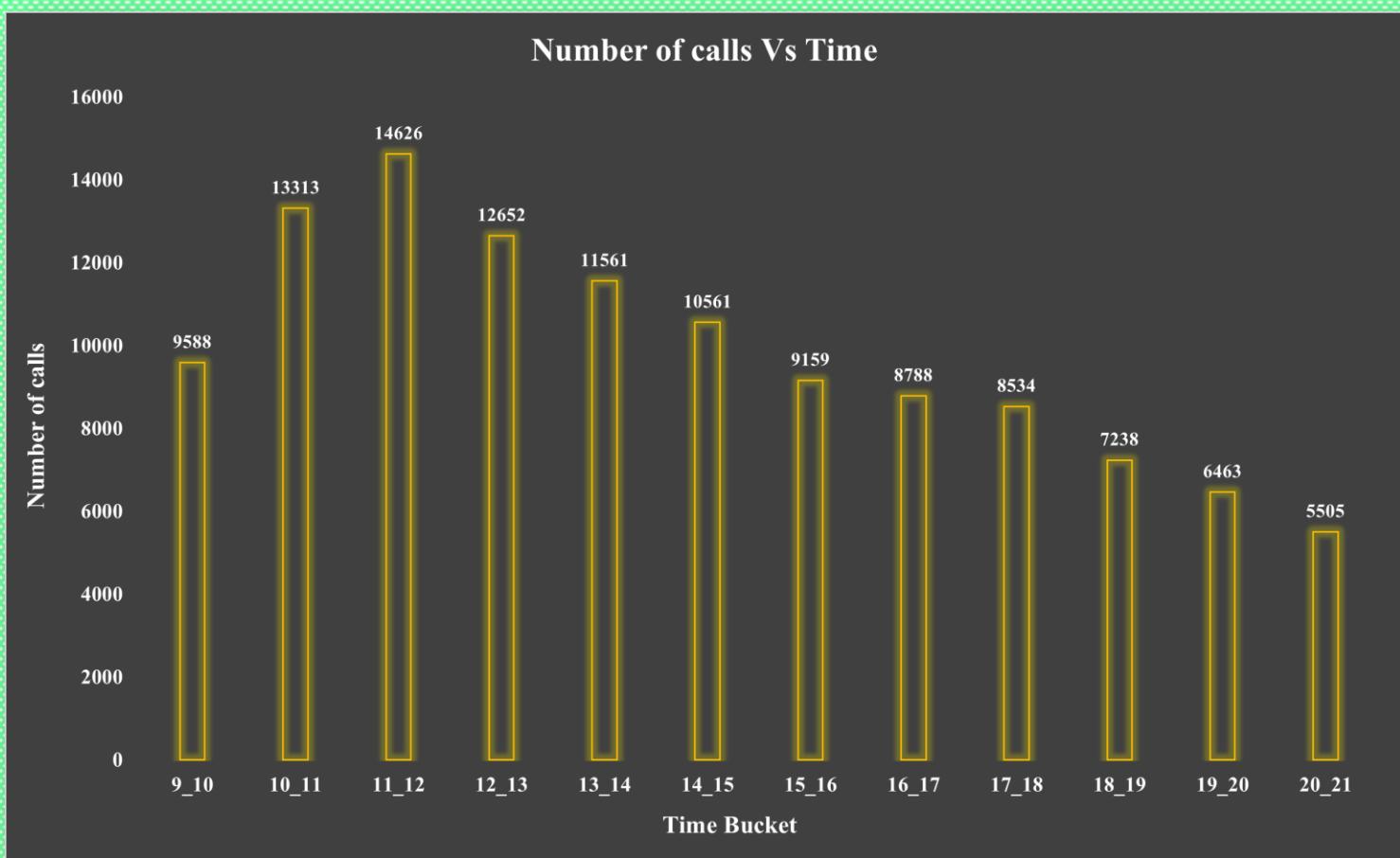


2. Call Volume Analysis:

Can you create a chart or graph that shows the number of calls received in each time bucket?

From the bar chart we can infer that time bucket 11_12 i.e. 11AM to 12PM has the highest count for total number incoming calls i.e. **14626**

| Time bucket | Number of calls |
|-------------|-----------------|
| 9_10 | 9588 |
| 10_11 | 13313 |
| 11_12 | 14626 |
| 12_13 | 12652 |
| 13_14 | 11561 |
| 14_15 | 10561 |
| 15_16 | 9159 |
| 16_17 | 8788 |
| 17_18 | 8534 |
| 18_19 | 7238 |
| 19_20 | 6463 |
| 20_21 | 5505 |



3. Manpower Planning:

Your Task: What is the minimum number of agents required in each time bucket to reduce the abandon rate to 10%?

Assumptions:

| | | |
|-----------------------------------|------------|------------|
| Agents working Hrs | 9 | Hrs |
| Break | 1.5 | Hrs |
| Agents working Hrs | 7.5 | Hrs |
| Agents Time spent on calls | 4.5 | Hrs |

Insights:

| | | |
|----------------------------------------------------|-------------|----------------|
| Avg Call Volume per day(9:00 AM to 9:00 PM) | 5130 | Nos |
| Average call duration (9:00 AM to 9:00 PM) | 199 | Seconds |
| Total call duration for 90% Calls | 255 | Hrs |
| Agents required per day | 57 | Persons |

| Time Bucket | Count of Time | Percentage of calls inbound | Agents Required |
|--------------------|---------------|-----------------------------|-----------------|
| 9_10 | 9588 | 8% | 5 |
| 10_11 | 13313 | 11% | 6 |
| 11_12 | 14626 | 12% | 7 |
| 12_13 | 12652 | 11% | 6 |
| 13_14 | 11561 | 10% | 6 |
| 14_15 | 10561 | 9% | 5 |
| 15_16 | 9159 | 8% | 4 |
| 16_17 | 8788 | 7% | 4 |
| 17_18 | 8534 | 7% | 4 |
| 18_19 | 7238 | 6% | 3 |
| 19_20 | 6463 | 5% | 3 |
| 20_21 | 5505 | 5% | 3 |
| Grand Total | 117988 | 100% | 57 |



The minimum number of agents required to reduce the abandon rate to 10% is **57 agents**.

4.Night Shift Manpower Planning:

Propose a manpower plan for each time bucket throughout the day, keeping the maximum abandon rate at 10%.

Assumptions:

| | | |
|-----------------------------------|------------|------------|
| Agents working Hrs | 9 | Hrs |
| Break | 1.5 | Hrs |
| Agents working Hrs | 7.5 | Hrs |
| Agents Time spent on calls | 4.5 | Hrs |

Insights:

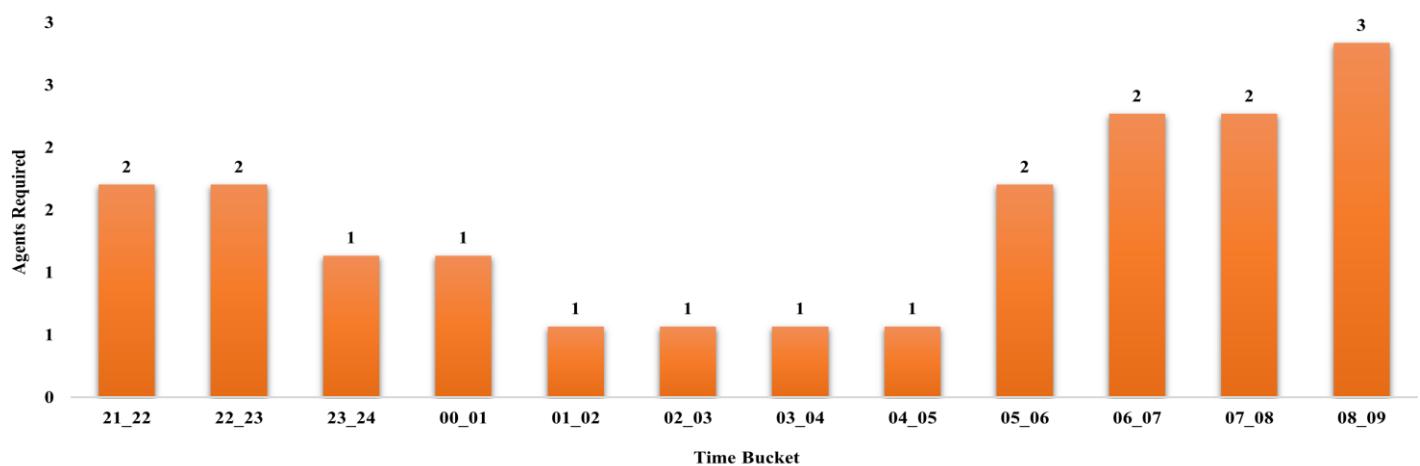
| | | |
|---------------------------------------------|------|---------|
| Avg Call Volume per day(9:00 AM to 9:00 PM) | 5130 | Nos |
| Average call duration (9:00 AM to 9:00 PM) | 199 | Seconds |
| Total call duration for 90% Calls | 255 | Hrs |
| Agents required per day | 57 | Persons |

Agents required (9:00 PM to 9:00 AM):

| | | |
|-----------------------------------------|------|---------|
| Avg Call Volume (9:00 PM to 9:00 AM) | 1539 | Nos |
| Total Call duration (9:00 PM to 9:00AM) | 77 | Hours |
| Agents required (9:00 PM to 9:00AM) | 17 | Persons |

| Time_Bucket | Calls in nights | Percentage of calls | Required agents in night |
|-------------|-----------------|---------------------|--------------------------|
| 21_22 | 3 | 10% | 2 |
| 22_23 | 3 | 10% | 2 |
| 23_24 | 2 | 7% | 1 |
| 00_01 | 2 | 7% | 1 |
| 01_02 | 1 | 3% | 1 |
| 02_03 | 1 | 3% | 1 |
| 03_04 | 1 | 3% | 1 |
| 04_05 | 1 | 3% | 1 |
| 05_06 | 3 | 10% | 2 |
| 06_07 | 4 | 13% | 2 |
| 07_08 | 4 | 13% | 2 |
| 08_09 | 5 | 17% | 3 |
| Total | 30 | 100% | 17 |

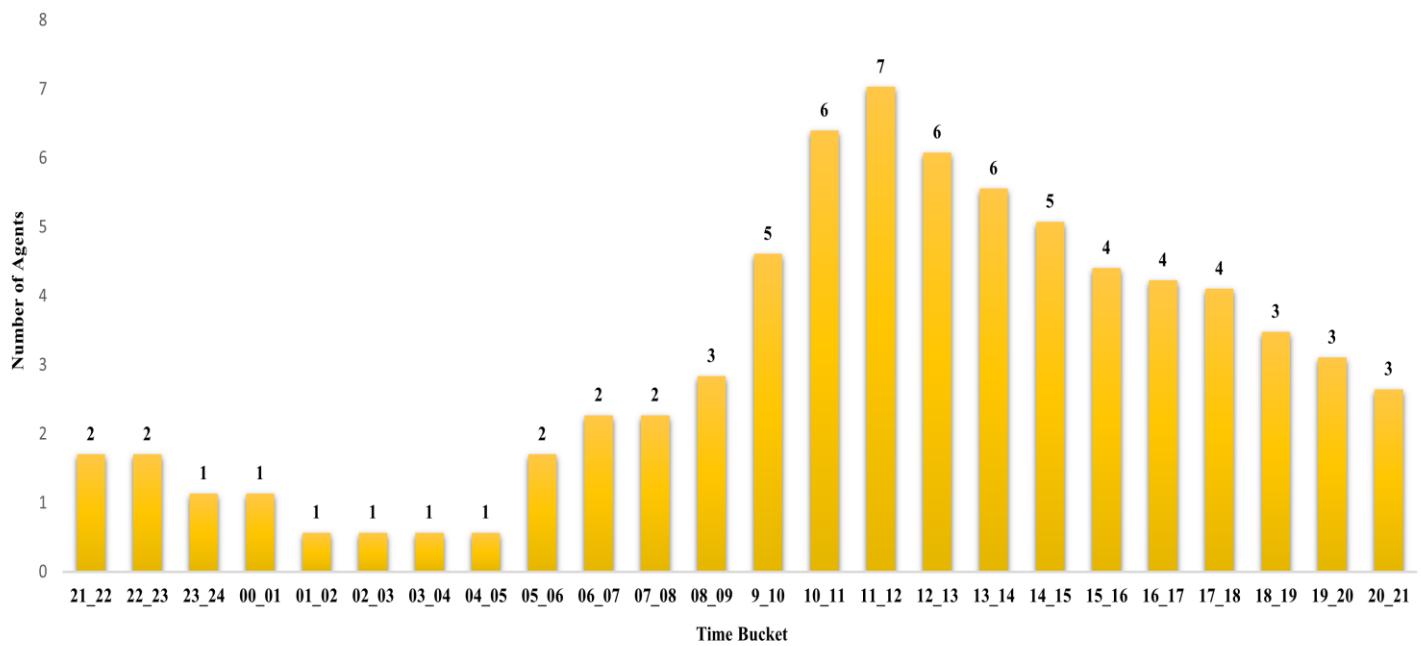
Required agents in Night



Therefore the ABC company needs **17 Agents** to answer the consumer calls at night time by keeping the answered rate to 90% / Abandon rate to 10%.

| Time_Bucket | Required Agents throughout the day |
|-------------|------------------------------------|
| 21_22 | 2 |
| 22_23 | 2 |
| 23_24 | 1 |
| 00_01 | 1 |
| 01_02 | 1 |
| 02_03 | 1 |
| 03_04 | 1 |
| 04_05 | 1 |
| 05_06 | 2 |
| 06_07 | 2 |
| 07_08 | 2 |
| 08_09 | 3 |
| 9_10 | 5 |
| 10_11 | 6 |
| 11_12 | 7 |
| 12_13 | 6 |
| 13_14 | 6 |
| 14_15 | 5 |
| 15_16 | 4 |
| 16_17 | 4 |
| 17_18 | 4 |
| 18_19 | 3 |
| 19_20 | 3 |
| 20_21 | 3 |

Required Agents throughout the day



Therefore the ABC company needs **74 Agents per day** to answer the consumer calls for day as well as the night time keeping the answered rate to 90% / Abandon rate to 10%.

CONCLUSION

From the above analysis, we can infer that the provided dataset has several meaningful insights which can help the ABC company to optimize customer service department and it can help improve the customer satisfaction.

Link for all the projects:

<https://drive.google.com/drive/folders/1HNYuAyzY5zMzVKOtDuT8rvizwe4oov8d?usp=sharing>

Link for Portfolio:

<https://drive.google.com/drive/folders/1MyoBEz7u9Pel8-hSEEluqP3uirO9dNKp?usp=sharing>