

# DATA ANALYSIS PORTFOLIO

SOWMIYA R



## **PROFESSIONAL BACKGROUND**

**Currently in final year pursuing B.Tech. Information Technology. I have secured 8.78 CGPA (till 6th Sem) and have several skills including Data Analysis, Python, MySQL, Excel. I have worked on various personal projects related to Data Analysis, Excel and Python.**

**Also, I have participated in various challenges and competitions on Hacker rank. Also, I have worked on Department Student Record, Counselling and Mentor system as part of mini project in Semester 5. Currently, I am working on implementation of mini project which is a group project for the Final year.**

**As I am a fresher it would be great to experience the real challenges of the corporate world and understand how things work. Being a fresher, I think I am very flexible and adaptive to learn new things. I'm excited to embark on a career that allows me to apply my knowledge and skills in a meaningful way. In my free time. And I believe by putting significant efforts, I will learn.**

# TABLE OF CONTENTS

## Catalog

Professional Background.....	2
Table Of Contents.....	3-4

### Data Analytics Process

Description.....	5
Design.....	6
Conclusion.....	8

### Instagram User Analytics

Description.....	10
The Problem.....	11
Design.....	12
Findings.....	13-21
Conclusion.....	21

### Operation Analytics and Investigating Metric Spike

Description.....	23
The Problem.....	24
Design.....	25
Findings.....	26-35
Conclusion.....	36

### Hiring Process Analytics

Description.....	38
The Problem.....	39-40
Design.....	41
Findings.....	42-49
Conclusion.....	50

Click here to enter text.

## **IMDB Movie Analysis**

Description.....	52
The Problem.....	53
Design.....	54
Findings.....	55-80
Conclusion.....	81

## **Bank Loan Case Study**

Description.....	83
The Problem.....	83
Design.....	84
Findings.....	86-96
Conclusion.....	96-98

## **Impact Of Car Features**

Description.....	100
The Problem.....	101-102
Design.....	103
Findings.....	104-118
Conclusion.....	118

## **ABC Call Volume Trend**

Description.....	120
The Problem.....	121
Design.....	122
Findings.....	122-127
Conclusion.....	128

<b>Appendices.....</b>	<b>129-131</b>
------------------------	----------------

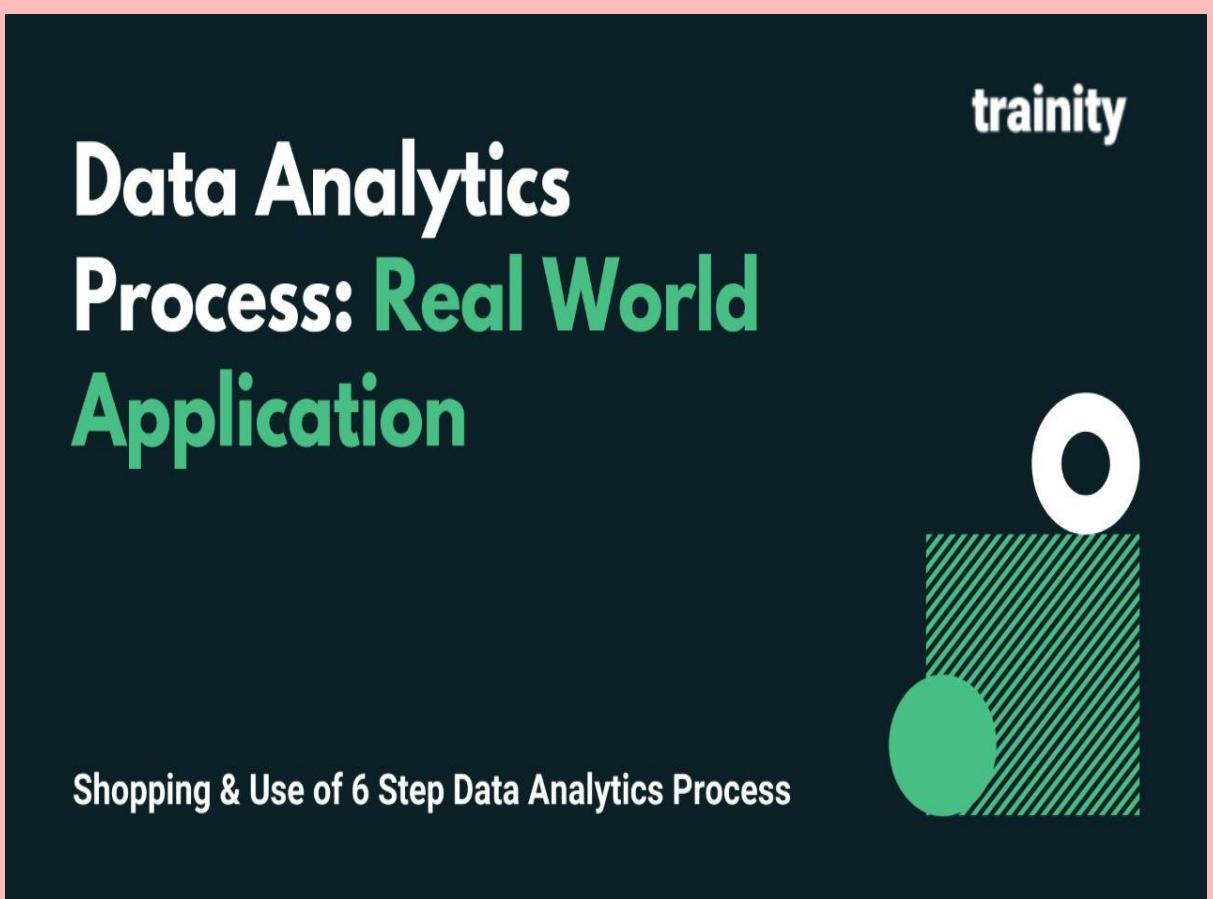


## 1. DATA ANALYTICS PROCESS

### DESCRIPTION

We use Data Analytics in everyday life without even knowing it.

For eg : Going to a market to buy something .



A slide titled "Data Analytics Process: Real World Application" from "trainity". The title is in large white and green text. Below the title, there is a subtitle "Shopping & Use of 6 Step Data Analytics Process" in white. To the right of the text, there is a graphic element consisting of a white circle and a green square with diagonal stripes.

## **DESIGN**

Tour package suppliers are continually looking for ways to improve customer experiences, optimize services, and keep ahead of market trends in an increasingly competitive travel business. By employing data-driven insights to make educated decisions, data analytics has emerged as a valuable tool for achieving these aims.

## **HERE THE PROCESS ARE**

### **Planning:**

- Analyze the needs of the customer, including where they are satisfied, costeffective, for lodging, sightseeing, and food.
- Where is it safe to travel?
- Which package is both affordable and fun?

### **Prepare:**

Collecting the information linked to customer requirement to serve the purpose :

- By searching the good location
- By finding the good accomodation
- By searching the transportation and food

Which are affordable to customer.

## **Process:**

By using these information I need to process all of this information and make judgments after receiving it.

This includes by comparing the:

- cost and travel distance
  - Quality of service
  - Quantity of service
  - Safety measures

So these are things used in the process method.

## Analyze:

Here we are analyzing the data by comparing the place by using the data table:

**PLACE COST TRAVEL TIME DISTANCE NUMBER  
OF SIGHTSEEING**

PLACE

1. From chennai to Wayanad 3500 11hr 16min 614km 8
  2. From chennai to Alapuzha 4000 14hr 33min 736km 10
  3. From chennai to goa 4500 16hr 36min 911km 9
  4. From chennai to bangalore 3000 7hr 13min 346km 5
  5. From chennai to manali 8000 48hr 2687km 13

### **Share:**

By sharing the analyzed data to customer to choose their tour package and also need to get feedback.

### **Act:**

Now they can choose their tour package which they can afford easily with the good quality of a service and ending with the payment for the tour package by using UPI option, g pay or any other online payment option.

### **Conclusion:**

- Finally customer satisfies with good quality of service.
- They found the affordable and enjoyable quantity of service.
- They found satisfied tour package.

### **Data Analytics Steps**

- i. Plan
- ii. Prepare
- iii. Process
- iv. Analyze
- v. Share
- vi. Act

# INSTAGRAM USER ANALYTICS



## **DESCRIPTION**

Imagine you're a data analyst working with the product team at Instagram. Your role involves analyzing user interactions and engagement with the Instagram app to provide valuable insights that can help the business grow.

User analysis involves tracking how users engage with a digital product, such as a software application or a mobile app. The insights derived from this analysis can be used by various teams within the business. For example, the marketing team might use these insights to launch a new campaign, the product team might use them to decide on new features to build, and the development team might use them to improve the overall user experience.

In this project, you'll be using SQL and MySQL Workbench as your tool to analyze Instagram user data and answer questions posed by the management team. Your insights will help the product manager and the rest of the team make informed decisions about the future direction of the Instagram app.

Remember, the goal of this project is to use your SQL skills to extract meaningful insights from the data. Your findings could potentially influence the future development of one of the world's most popular social media platforms.

## The Problem

### A) Marketing Analysis:

1. **Rewarding Most Loyal Users:** Finding the 5 oldest users of the Instagram from the database provided with the joined date.
2. **Remind Inactive Users to Start Posting:** By sending them promotional emails to post their 1st photo. Find the users who have never posted a single photo on Instagram.
3. **Declaring Contest Winner:** The team started a contest and the user who gets the most likes on a single photo will win the contest now they wish to declare the winner. Identify the winner of the contest and provide their details to the team.
4. **Hash tag Researching:** A partner brand wants to know, which hash tags to use in the post to reach the most people on the platform. Identify and suggest the top 5 most commonly used hash tags on the platform.
5. **Launch AD Campaign:** The team wants to know, which day would be the best day to launch ADs. What day of the week

do most users register on? Provide insights on when to schedule an ad campaign.

## **B) Investor Metrics:**

1. **User Engagement:** Are users still as active and post on Instagram or they are making fewer posts. Provide how many times does average user posts on Instagram. Also, provide the total number of photos on Instagram/total number of users.
2. **Bots & Fake Accounts:** The investors want to know if the platform is crowded with fake and dummy accounts. Provide data on users (bots) who have liked every single photo on the site (since any normal user would not be able to do this).

## **DESIGN**

- Using the 'create db' function of MySQL create a database
- Then add tables and column names
- Then add the values into them using the 'insert into' function of MySQL
- By using the 'select' command we can query the desired output

--> **MySQL Workbench 8.0 CE**

## Findings 1

**1.Rewarding Most Loyal Users:** Finding the 5 oldest users of the Instagram from the database provided with the joined date.

```
SELECT *FROM users  
ORDER BY created_at DESC LIMIT 5;
```

```
/*Table: users*/
CREATE TABLE users(
    id INT NOT NULL AUTO_INCREMENT PRIMARY KEY,
    tag_name VARCHAR(255) UNIQUE NOT NULL,
    created_at TIMESTAMP DEFAULT NOW()
);

/*Table: photos - Tags*/
CREATE TABLE photo_tags(
    photo_id INT NOT NULL,
    tag_id INT NOT NULL,
    FOREIGN KEY(photo_id) REFERENCES photos(id),
    FOREIGN KEY(tag_id) REFERENCES tag(id),
    PRIMARY KEY(photo_id,tag_id)
);

INSERT INTO users(username, created_at) VALUES ('Kenton_Kirlin', '2017-02-16 18:22:10.846'), ('Andrea_Purdy85', '2017-04-07 11:11:11'), ('Leland_Hettinger', '2017-04-08 03:03:03'), ('Audrey_Miller', '2017-04-09 04:04:04'), ('Carmen_Hayes', '2017-04-10 05:05:05');

INSERT INTO photos(image_url, user_id) VALUES ('http://ellijah.biz', 1), ('https://shanon.org', 1), ('http://vicky.biz', 1);

INSERT INTO follows(follower_id, followee_id) VALUES (2, 1), (2, 3), (2, 4), (2, 5), (2, 6), (2, 7), (2, 8), (2, 9), (2, 10);

INSERT INTO comments(comment_text, user_id, photo_id) VALUES ('unde at dolore', 2, 1), ('que es ducimus', 3, 1), ('alias', 4, 1);

INSERT INTO likes(user_id,photo_id) VALUES (2, 1), (5, 1), (9, 1), (10, 1), (11, 1), (14, 1), (19, 1), (21, 1), (24, 1), (25, 1);

INSERT INTO tag(tag_name) VALUES ('sunset'), ('photography'), ('sunrise'), ('landscape'), ('food'), ('foodie'), ('delicio');

INSERT INTO photo_tags(photo_id, tag_id) VALUES (1, 10), (1, 17), (1, 21), (1, 13), (1, 19), (2, 4), (2, 3), (2, 20), (2, 18);

SELECT *FROM users ORDER BY created_at DESC LIMIT 5;
```

Oldest 5 users joining date with their username and user id:

<b>id</b>	<b>username</b>	<b>created_at</b>
38	Jordyn.Jacobson2	2016-05-14 07:56:26
63	Elenor88	2016-05-08 01:30:41
67	Emilio_Bernier52	2016-05-06 13:04:30
80	Darby_Herzog	2016-05-06 00:14:21
95	Nicole71	2016-05-06 00:14:21

## Findings 2

**2.Remind Inactive Users to Start Posting:** By sending them promotional emails to post their 1st photo. Find the users who have never posted a single photo on Instagram.

```
select id,username  
from users  
WHERE id NOT IN (SELECT user_id FROM photos);
```

The screenshot shows a MySQL query being run on a database named 'ig\_clone'. The query is:

```
select id,username  
from users  
WHERE id NOT IN (SELECT user_id FROM photos);
```

The output of the query is a list of 26 user IDs and their corresponding usernames. The output table is as follows:

id	username
5	Aniya_Hackett
7	Kasandra_Homenick
14	Jaclyn81
21	Rocio33
24	Maxwell.Halvorson
25	Tierra.Trantow
34	Pearl7
36	Ollie_Ledner37
41	Mckenna17
45	David.Osinski47
49	Morgan.Kassulke
53	Linnea59
54	Duane60
57	Julien_Schmidt
66	Mike.Auer39
68	Franco_Keebler64
71	Nia_Haag
74	Hulda.Macejkovic
75	Leslie67

There are total 26 users who never posted a photo on the platform:

id	username
5	Aniya_Hackett
7	Kasandra_Homenick
14	Jaclyn81
21	Rocio33
24	Maxwell.Halvorson
25	Tierra.Trantow
34	Pearl7
36	Ollie_Ledner37
41	Mckenna17
45	David.Osinski47
49	Morgan.Kassulke
53	Linnea59
54	Duane60
57	Julien_Schmidt
66	Mike.Auer39
68	Franco_Keebler64
71	Nia_Haag
74	Hulda.Macejkovic
75	Leslie67

54	Duane60
57	Julien_Schmidt
66	Mike.Auer39
68	Franco_Keebler64
71	Nia_Haag
74	Hulda.Macejkovic
75	Leslie67
76	Janelle.Nikolaus81
80	Darby_Herzog
81	Esther.Zulauf61
83	Bartholome.Bernhard
89	Jessyca_West
90	Esmeralda.Mraz57
91	Bethany20

## Findings 3

- 3. Declaring Contest Winner:** The team started a contest and the user who gets the most likes on a single photo will win the contest now they wish to declare the winner. Identify the winner of the contest and provide their details to the team.

```

SELECT username, photos.id, photos.image_url, count(*) as total
FROM photos inner join likes on likes.photo_id = photos.id
inner join users on photos.user_id = users.id
GROUP BY photos.id
ORDER BY total DESC LIMIT 1;

```

```

47 /*Tags*/
48 CREATE TABLE tags(
49   id INTEGER AUTO_INCREMENT PRIMARY KEY,
50   tag_name VARCHAR(255) UNIQUE NOT NULL,
51   created_at TIMESTAMP DEFAULT NOW()
52 );
53
54
55 /*junction table: Photos - Tags*/
56 CREATE TABLE photo_tags(
57   photo_id INT NOT NULL,
58   tag_id INT NOT NULL,
59   FOREIGN KEY(photo_id) REFERENCES photos(id),
60   FOREIGN KEY(tag_id) REFERENCES tags(id),
61   PRIMARY KEY(photo_id,tag_id)
62 );
63
64
65
66 INSERT INTO users (username, created_at) VALUES ('Kenton_Kirlin', '2017-02-16 18:22:10.846'), ('Andre_Purdy85', '2017-04-07 10:45:00.000')
67
68 INSERT INTO photos(image_url, user_id) VALUES ('http://elijah.biz', 1), ('https://shanon.org', 1), ('http://vicky.biz', 1)
69
70 INSERT INTO follows(follower_id, followee_id) VALUES (2, 1), (2, 3), (2, 4), (2, 5), (2, 6), (2, 7), (2, 8), (2, 9), (2, 10)
71
72
73
74
75 INSERT INTO comments(comment_text, user_id, photo_id) VALUES ('unde at dolorem', 2, 1), ('quea ea ducimus', 3, 1), ('alias autem', 4, 1)
76
77
78 INSERT INTO likes(user_id,photo_id) VALUES (2, 1), (5, 1), (9, 1), (10, 1), (11, 1), (14, 1), (19, 1), (21, 1), (24, 1), (25, 1)
79
80
81 INSERT INTO tags(tag_name) VALUES ('sunset'), ('photography'), ('sunrise'), ('landscape'), ('food'), ('foodie'), ('delicious'), ('travel')
82
83
84 INSERT INTO photo_tags(photo_id, tag_id) VALUES (1, 18), (1, 17), (1, 21), (1, 13), (1, 19), (2, 4), (2, 3), (2, 20), (2, 1)
85
86 select username,photos.id,photos.image_url,count(*) as total from photos
87 inner join likes on likes.photo_id = photos.id
88 inner join users on photos.user_id = users.id group by photos.id order by total DESC limit 1;
89

```

User with ID: **Zack\_kemmer93** has won the contest with 48 likes for a single photo he had posted.

username	id	image_url	total
Zack_Kemmer93	145	https://jarret.name	48

## Findings 4

**4. Hash tag Researching:** A partner brand wants to know, which hash tags to use in the post to reach the most people on the platform. Identify and suggest the top 5 most commonly used hash tags on the platform.

```

SELECT tags.tag_name, count(*) as total FROM photo_tags
join tags on photo_tags.tag_id=tags.id
GROUP BY tags.id
ORDER BY total ASC LIMIT 6;

```

```

48 /*Tags*/
49 CREATE TABLE tags(
50   id INTEGER AUTO_INCREMENT PRIMARY KEY,
51   tag_name VARCHAR(255) UNIQUE NOT NULL,
52   created_at TIMESTAMP DEFAULT NOW()
53 );
54
55 /*junction table: Photos - Tags*/
56 CREATE TABLE photo_tags(
57   photo_id INT NOT NULL,
58   tag_id INT NOT NULL,
59   FOREIGN KEY(photo_id) REFERENCES photos(id),
60   FOREIGN KEY(tag_id) REFERENCES tags(id),
61   PRIMARY KEY(photo_id,tag_id)
62 );
63
64
65 INSERT INTO users (username, created_at) VALUES ('Kenton_Kirlin', '2017-02-16 18:22:10.846'), ('Andre_Purdy85', '2017-04-0
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91

```

Output:

tag_name	total
foodie	11
delicious	15
photography	16
stunning	16
sunrise	17
landscape	17

tag_name	total
foodie	11
delicious	15
photography	16
stunning	16
sunrise	17
landscape	17

## Findings 5

5. **Launch AD Campaign:** The team wants to know, which day would be the best day to launch ADs. What day of the week do most users register on? Provide insights on when to schedule an ad campaign.

```

SELECT day_week, count(day_week) AS date_week
FROM (SELECT dayname(created_at) AS day_week
      FROM users) as table1
GROUP BY day_week
ORDER BY date_week DESC LIMIT 7;

```

The screenshot shows a MySQL query editor interface. On the left, the code for the query is displayed. On the right, there are tabs for 'NEW', 'MySQL', and 'RUN'. Below the tabs, there is an 'STDIN' input field labeled 'Input for the program (Optional)'. The 'Output' section displays the results of the query:

day_week	date_week
Thursday	16
Sunday	16
Friday	15
Tuesday	14
Monday	14
Wednesday	13
Saturday	12

day_week	date_week
Thursday	16
Sunday	16
Friday	15
Tuesday	14
Monday	14
Wednesday	13
Saturday	12

Thursdays and Sundays are two days with most user signup on Instagram.



## Findings 6

1. **User Engagement:** Are users still as active and post on Instagram or they are making fewer posts. Provide how many times does average user posts on Instagram. Also, provide the total number of photos on Instagram/total number of users.

**Total number of photos on Instagram / Total number of users**

```
SELECT (SELECT count(*) FROM photos)/(SELECT count(*) FROM users) AS average;
```

The screenshot shows a web-based MySQL query editor. The code area contains the provided SQL query. The results section shows the output of the query: 'avg' followed by '2.5700'. The interface includes tabs for 'queries.sql' and '+', and buttons for 'NEW', 'MYSQL', 'RUN', and 'STDIN'.

```
47
48 /*Tags*/
49 CREATE TABLE tags(
50   id INTEGER AUTO_INCREMENT PRIMARY KEY,
51   tag_name VARCHAR(255) UNIQUE NOT NULL,
52   created_at TIMESTAMP DEFAULT NOW()
53 );
54
55 /*junction table: Photos -> Tags*/
56 CREATE TABLE photo_tags(
57   photo_id INT NOT NULL,
58   tag_id INT NOT NULL,
59   FOREIGN KEY(photo_id) REFERENCES photos(id),
60   FOREIGN KEY(tag_id) REFERENCES tags(id),
61   PRIMARY KEY(photo_id,tag_id)
62 );
63
64
65 INSERT INTO users (username, created_at) VALUES ('Kenton_Kirlin', '2017-02-16 18:22:10.846'), ('Andre_Purdy85', '2017-04-0
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
```

```
ig_clone
```

Output:  
avg  
2.5700

average
2.5700

## Findings 7

**2. Bots & Fake Accounts:** The investors want to know if the platform is crowded with fake and dummy accounts. Provide data on users (bots) who have liked every single photo on the site (since any normal user would not be able to do this).

```
SELECT user_id, username, count(*) AS count_likes FROM users
inner join likes on users.id = likes.user_id
GROUP BY likes.user_id
HAVING count_likes = (SELECT count(*) FROM photos) ;
```

The screenshot shows the execution of the provided SQL query in an online MySQL compiler. The results are as follows:

user_id	username	count_likes
5	Aniya_Hackett	257
14	Jaclyn81	257
21	Rocio33	257
24	Maxwell.Halvorson	257

user_id	username	count_likes
5	Aniya_Hackett	257
14	Jaclyn81	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

## CONCLUSION

I would like to conclude that not only Instagram but many other social media and commercial firms use such Analysis to find the insights from their customer data which in turn help the firms to find the customers who will be an Asset to the firm in the future and not some Liability.

Such Analysis and sorting of the customer base is done at an weekly, monthly, quarterly or yearly basis as per the needs of the business firms so as to maximize their profits in future with minimal cost to the company.

**Instagram**



# Operation Analytics and Investigating Metric Spike

(Advanced sql)

**trainity**

## Operation Analytics & Investigating metric spike case study

**Analysis**

Category	Value
a	1
b	5
c	3
d	10

Legend: Marketing (orange), Developers (cyan), HR (red), Design (purple)

**Metric Spike**

Legend: Available (red), Target (blue)

Date	Available	Target
Oct 2021	6	4
Nov 2021	7	5
Dec 2021	6	5
Jan 2022	3	3
Feb 2022	7	6
Mar 2022	8	6

**Employees**

Aug 25-Sept 25 ✓

Category	Value
Inactive	254
Active	3000
Total	3254

## **DESCRIPTION**

Operational Analytics is a crucial process that involves analyzing a company's end-to-end operations. This analysis helps identify areas for improvement within the company. As a Data Analyst, you'll work closely with various teams, such as operations, support, and marketing, helping them derive valuable insights from the data they collect.

One of the key aspects of Operational Analytics is investigating metric spikes. This involves understanding and explaining sudden changes in key metrics, such as a dip in daily user engagement or a drop in sales. As a Data Analyst, you'll need to answer these questions daily, making it crucial to understand how to investigate these metric spikes.

In this project, you'll take on the role of a Lead Data Analyst at a company like Microsoft. You'll be provided with various datasets and tables, and your task will be to derive insights from this data to answer questions posed by different departments within the company. Your goal is to use your advanced SQL skills to analyze the data and provide valuable insights that can help improve the company's operations and understand sudden changes in key metrics.

## The Problem

### Case Study 1 (Job Data)

- **Number of jobs reviewed:** Amount of jobs reviewed over time.
- **Throughput:** It is the no. of events happening per second.
- **Percentage share of each language:** Share of each language for different contents.
- **Duplicate rows:** Rows that have the same value present in them.

### Case Study 2 (Investigating metric spike)

- **User Engagement:** To measure the activeness of a user. Measuring if the user finds quality in a product/service.
- **User Growth:** Amount of users growing over time for a product.
- **Weekly Retention:** Users getting retained weekly after signing-up for a product.
- **Weekly Engagement:** To measure the activeness of a user. Measuring if the user finds quality in a product/service weekly.
- **Email Engagement:** Users engaging with the email service.

## **Design**

- Using the 'create db' function of MySQL create a database
- Then add tables and column names
- Then add the values into them using the 'insert into' function of MySQL
- By using the 'select' command we can query the desired output

### **Steps taken to load the data into the data base**

Software used for querying the results

--> **MySQL Workbench 8.0 CE**

Software used for analyzing using Bar plots

--> **Microsoft Excel**

## Findings 1

1. **Number of jobs reviewed:** Amount of jobs reviewed over time.

Our task is to Calculate the number of jobs reviewed per hour per day for November 2020.

```
SELECT ds AS Dates,  
(round(count(job_id)/sum(time_spent))*3600) as 'Jobs reviewed  
per hr'  
from `sql project-1 table (1)`  
where ds between '2020-11-01' and '2020-11-30'  
group by ds;
```

	Dates	Jobs reviewed per hr
▶	2020-11-30	0
^	2020-11-29	0
	2020-11-28	0
	2020-11-27	0
	2020-11-26	0
	2020-11-25	0

## Findings 2

**2.Throughput:** It is the no. of events happening per second.

Our task is to Calculate 7 day rolling average of throughput? For throughput, do you prefer daily metric or 7-day rolling and why?

```
SELECT (ROUND(COUNT(event)/sum(time_spent),2)) as  
Throughput_weekly  
FROM `sql project-1 table (1)`;
```

The weekly throughput is 0.03

Result Grid	
	Throughput_weekly
▶	0.03

```
select ds as dates, (ROUND(COUNT(event)/sum(time_spent),2))  
as Throughput_daily  
FROM `sql project-1 table (1)`  
group by ds  
order by ds;
```

Result Grid	
dates	Throughput_daily
▶ 2020-11-25	0.02
2020-11-26	0.02
2020-11-27	0.01
2020-11-28	0.06
2020-11-29	0.05
2020-11-30	0.05

## Findings 3

**3.Percentage share of each language:** Share of each language for different contents.

Our task is to calculate the percentage share of each language in the last 30 days.

```
SELECT 'language' AS Languages, (ROUND (100*  
COUNT(*)/total,2)) as percentage_of_lang  
  
FROM `sql project-1 table (1)`  
  
CROSS JOIN (SELECT COUNT(*) as total FROM `sql project-1  
table (1)` ) sub_total  
  
GROUP BY languages;
```

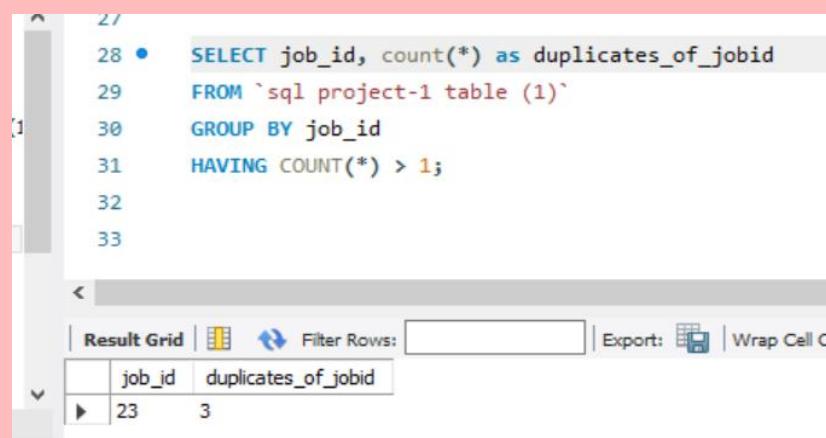
	Languages	Percentage
▶	English	12.50
	Arabic	12.50
	Persian	37.50
	Hindi	12.50
	French	12.50
	Italian	12.50

## Findings 4

**4.Duplicate rows:** Rows that have the same value present in them.

Our task is to display duplicates from the table.

```
SELECT job_id, count(*) as duplicates_of_jobid  
FROM `sql project-1 table (1)`  
GROUP BY job_id  
HAVING COUNT(*) > 1;
```



The screenshot shows a SQL editor interface. The code area contains the following SQL query:

```
27  
28 •  SELECT job_id, count(*) as duplicates_of_jobid  
29   FROM `sql project-1 table (1)`  
30   GROUP BY job_id  
31   HAVING COUNT(*) > 1;  
32  
33
```

Below the code, there is a toolbar with buttons for 'Result Grid' and 'Filter Rows'. The result grid displays the following data:

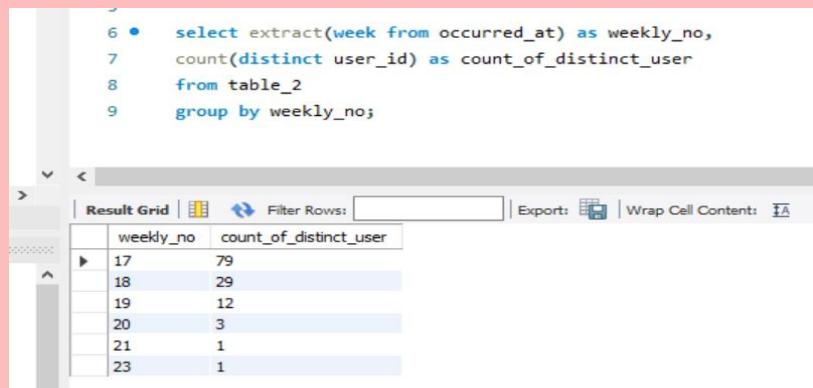
job_id	duplicates_of_jobid
23	3

## Findings 5

1. **User Engagement:** To measure the activeness of a user. Measuring if the user finds quality in a product/service.

Our task is to calculate the weekly user engagement.

```
SELECT EXTRACT(week FROM occurred_at) AS weekly_no,  
COUNT(distinct user_id) AS count_of_distinct_user  
FROM table_2  
GROUP BY weekly_no;
```



The screenshot shows a database query editor with the following content:

```
6 •  select extract(week from occurred_at) as weekly_no,  
7      count(distinct user_id) as count_of_distinct_user  
8      from table_2  
9      group by weekly_no;
```

Below the code is a result grid table:

weekly_no	count_of_distinct_user
17	79
18	29
19	12
20	3
21	1
23	1

## Findings 6

2. **User Growth:** Amount of users growing over time for a product.

Our task is to calculate the user growth for product.

```

SELECT year, num_of_week, num_of_active_users,
SUM(num_of_active_users) OVER(ORDER BY year,
num_of_week rows BETWEEN unbounded preceding and current
row)
AS cumm_of_active_users
FROM
(SELECT
EXTRACT(year FROM a.activated_at) AS year,
EXTRACT(week FROM a.activated_at) AS num_of_week,
COUNT(distinct user_id) AS num_of_active_users
FROM table_1 a
WHERE state='active'
GROUP BY year, num_of_week
ORDER BY year, num_of_week
)a;

```

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

	year	num_of_week	num_of_active_users	cumm_of_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	42	295
	2013	8	34	329
	2013	9	43	372

Result 3 ×

## Findings 7

**3. Weekly Retention:** Users getting retained weekly after signing-up for a product.

Our task is to calculate the weekly retention of users-sign up cohort.

```
select count(user_id),
       sum(case when retention_week = 1 then 1 else 0 end) as
       retention_per_week
from
(
  select (a.user_id,
          a.sign_up_week,
          b.engagement_week,
          b.engagement_week - a.sign_up_week) as
          retention_week
    from
    (
      (select distinct user_id, extract(week from occurred_at) as
       week_signup
       from table_2
      where event_type = 'signup_flow'
      and event_name = 'complete_signup'
      and extract(week from occurred_at)=18)a
     left join
      (select distinct user_id, extract(week from occurred_at) as
       engagement_week
       from table_2
      where event_type = 'engagement')b
     on a.user_id = b.user_id
    )
  )
group by user_id
order by user_id;
```

year	week	device	count(distinct user_id)
NONE	NONE	acer aspire desktop	198
NONE	NONE	acer aspire notebook	338
NONE	NONE	amazon fire phone	89
NONE	NONE	asus chromebook	355
NONE	NONE	dell inspiron desktop	360
NONE	NONE	dell inspiron notebook	677
NONE	NONE	hp pavilion desktop	339
NONE	NONE	htc one	196
NONE	NONE	ipad air	478
NONE	NONE	ipad mini	292
NONE	NONE	iphone 4s	409
NONE	NONE	iphone 5	1025
NONE	NONE	iphone 5s	626
NONE	NONE	kindle fire	205
NONE	NONE	lenovo thinkpad	1309
NONE	NONE	mac mini	150
NONE	NONE	macbook air	950
NONE	NONE	macbook pro	1952
NONE	NONE	nexus 10	273
NONE	NONE	nexus 5	621
NONE	NONE	nexus 7	355
NONE	NONE	nokia lumia 635	211
NONE	NONE	samsung galaxy tablet	107
NONE	NONE	samsung galaxy note	119
NONE	NONE	samsung galaxy s4	803
NONE	NONE	windows surface	182

## Findings 8

**4. Weekly Engagement:** To measure the activeness of a user. Measuring if the user finds quality in a product/service weekly.

Our task is to calculate the weekly engagement per device.

```
select extract(week from occurred_at)as week_signup, extract(year from occurred_at)as year, device,
count(distinct user_id)as 'count_of_user'
from table_2
where event_type='engagement'
```

group by 1,2,3

order by 1,2,3;

The screenshot shows a MySQL Workbench result grid with the following data:

	week_signup	year	device	count_of_user
▶	17	2014	acer aspire desktop	1
	17	2014	acer aspire notebook	2
	17	2014	amazon fire phone	1
	17	2014	asus chromebook	3
	17	2014	dell inspiron desktop	1
	17	2014	dell inspiron notebook	4
	17	2014	hp pavilion desktop	2
	17	2014	htc one	2
	17	2014	ipad air	1
	17	2014	ipad mini	3
	17	2014	iphone 4s	3
	17	2014	iphone 5	9
	17	2014	iphone 5s	5
	17	2014	lenovo thinkpad	6
	17	2014	mac mini	1

## Findings 9

**5.Email Engagement:** Users engaging with the email service.

Our task is to calculate the email engagement metrics.

select

100.0 \* sum(case when email\_cat = 'email\_opened' then 1 else 0 end)

/ sum(case when email\_cat = 'email\_sent' then 1 else 0 end)

as email\_open\_rate,

100.0 \* sum(case when email\_cat = 'email\_clicked' then 1 else 0 end)

```

/sum(case when email_cat = 'email_sent' then 1 else 0 end)

as email_click_rate

from

(
select *,

case when action in ('sent_weekly_digest',
'sent_reengagement_email')

then 'email_sent'

when action in ('email_open')

then 'email_opened'

when action in ('email_clickthrough')

then 'email_clicked'

end as email_cat

from email_events

)a;

```

The screenshot shows a database query result grid with the following data:

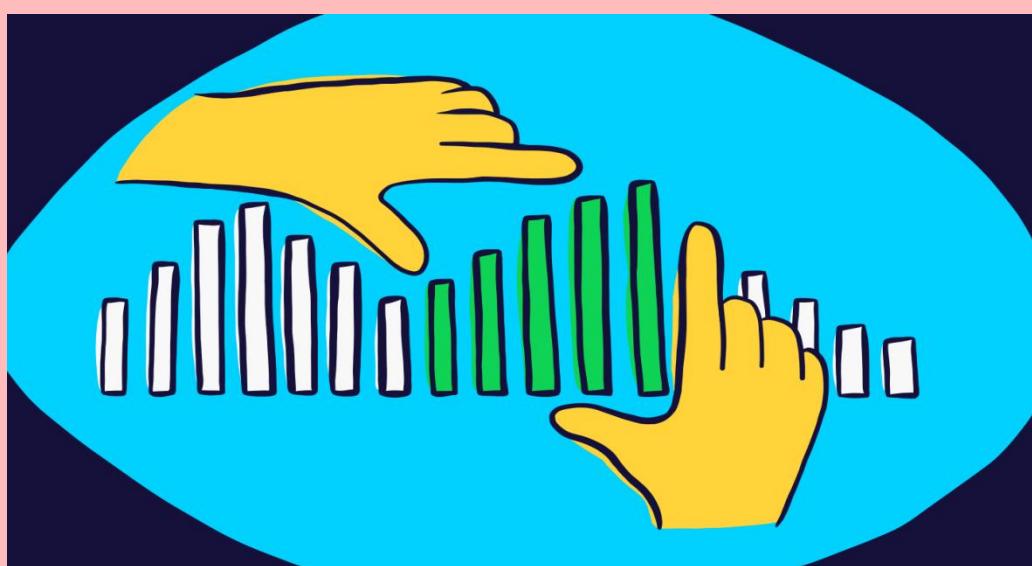
	email_open_rate	email_click_rate
▶	31.01294	10.99955

## Conclusion

I would like to conclude that Operation Analytics and Investigating Metric Spike are very necessary and they must be done on daily, weekly, Monthly, Quarterly or Yearly basis based on the Business needs of the firm.

Also, any firm/entity must focus on the Email Engagement with the customers; the firm must use catchy headings along with reasonable discounts and coupons so as to increase their existing customer base.

Also any firm must have a separate department(if possible) so as to hear out to the problems of those Visitors who had left the Sign-up Process in between, the firm must guide them so as to convert them from Visitors to Customers.





## Hiring Process Analytics



## Description

Hiring process is the fundamental and the most important function of a company. Here, the MNC's get to know about the major underlying trends about the hiring process. Trends such as- number of rejections, number of interviews, types of jobs, vacancies etc. are important for a company to analyse before hiring freshers or any other individual. Thus, making an opportunity for a Data Analyst job here too!

Being a Data Analyst, your job is to go through these trends and draw insights out of it for hiring department to work upon.

You are working for a MNC such as Google as a lead Data Analyst and the company has provided with the data records of their previous hiring and have asked you to answer certain questions making sense out of that data.

## The Problem

### Data Analytics Tasks:

After downloading the dataset, use Excel to answer the below questions:

**A. Hiring Analysis:** The hiring process involves bringing new individuals into the organization for various roles.

**Your Task:** Determine the gender distribution of hires. How many males and females have been hired by the company?

**B. Salary Analysis:** The average salary is calculated by adding up the salaries of a group of employees and then dividing the total by the number of employees.

**Your Task:** What is the average salary offered by this company? Use Excel functions to calculate this.

**C. Salary Distribution:** Class intervals represent ranges of values, in this case, salary ranges. The class interval is the difference between the upper and lower limits of a class.

**Your Task:** Create class intervals for the salaries in the company. This will help you understand the salary distribution.

**D. Departmental Analysis:** Visualizing data through charts and plots is a crucial part of data analysis.

**Your Task:** Use a pie chart, bar graph, or any other suitable visualization to show the proportion of people working in different departments.

**E. Position Tier Analysis:** Different positions within a company often have different tiers or levels.

**Your Task:** Use a chart or graph to represent the different position tiers within the company. This will help you understand the distribution of positions across different tiers.

## **Design**

Before starting the actual analysis I have:-

- Firstly I made a copy of the raw data where I can perform the Analysis so that what ever changes I made it will not affect the original data.
- Secondly I looked for blank spaces and NULL values if any.
- Then I had imputed the numerical blank and NULL cells with mean of the column(if no outliers existed for that particular column) or with median (if outliers existed for that column)
- Then I looked for if any outliers exists and replaced them with the median of the particular column where the outlier existed.
- Then for blank cells of categorical variables I had replaced with the variable with the highest count Then I looked for duplicate rows and removed them if any.
- Then I removed the irrelevant columns(data) from the dataset which was not necessary for doing the analysis.

Software used for doing the overall Analysis:-

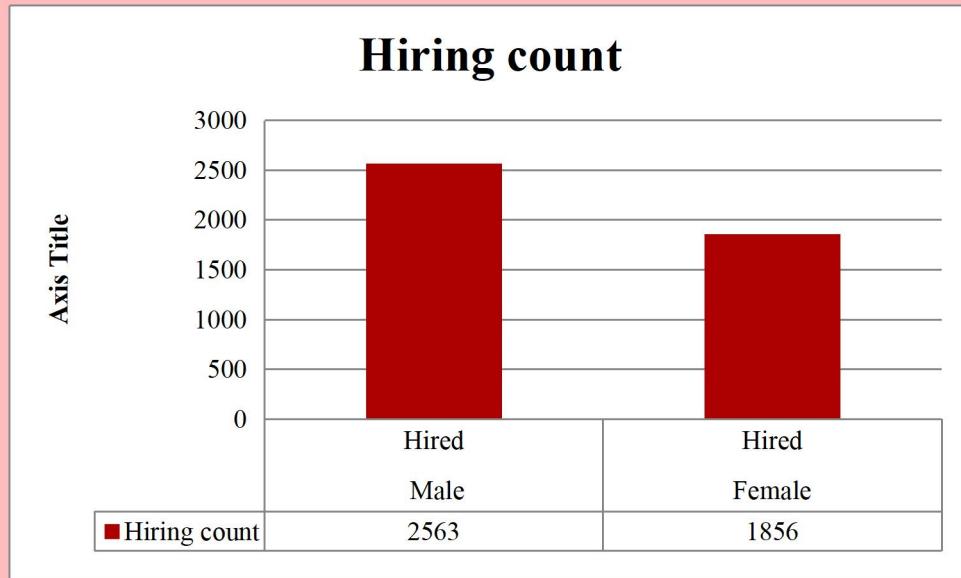
----> **Microsoft Excel**

## Findings 1

1. **Hiring:** Process of intaking of people into an organization for different kinds of positions.

**Our Task:** How many males and females are Hired ?

event_name	status	Hiring count
Male	Hired	2563
Female	Hired	1856



From the above bar plot , we can observe that there are 2563 males and 1856 females have been hired for different positions of organizations of the company.

## Findings 2

2. **Average Salary:** Adding all the salaries for a select group of employees and then dividing the sum by the number of employees in the group.

**Our Task :** What is the average salary offered in this company ?

Average salary	49983.03
Median of salary	49625
minimum salary	100
maximum salary	400000

From above , the average salary offered by this company was **49983.03**

## Findings 3

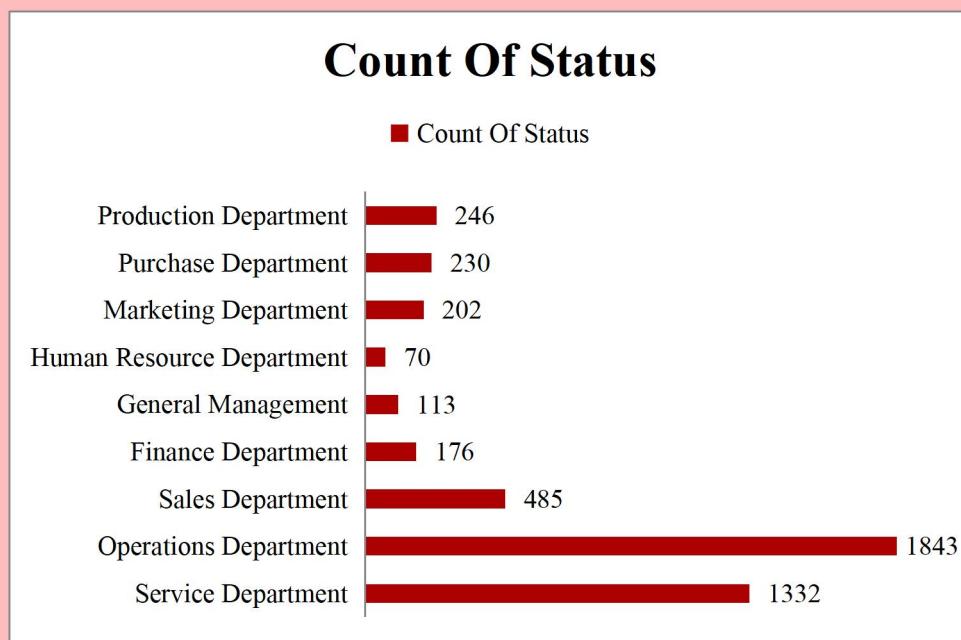
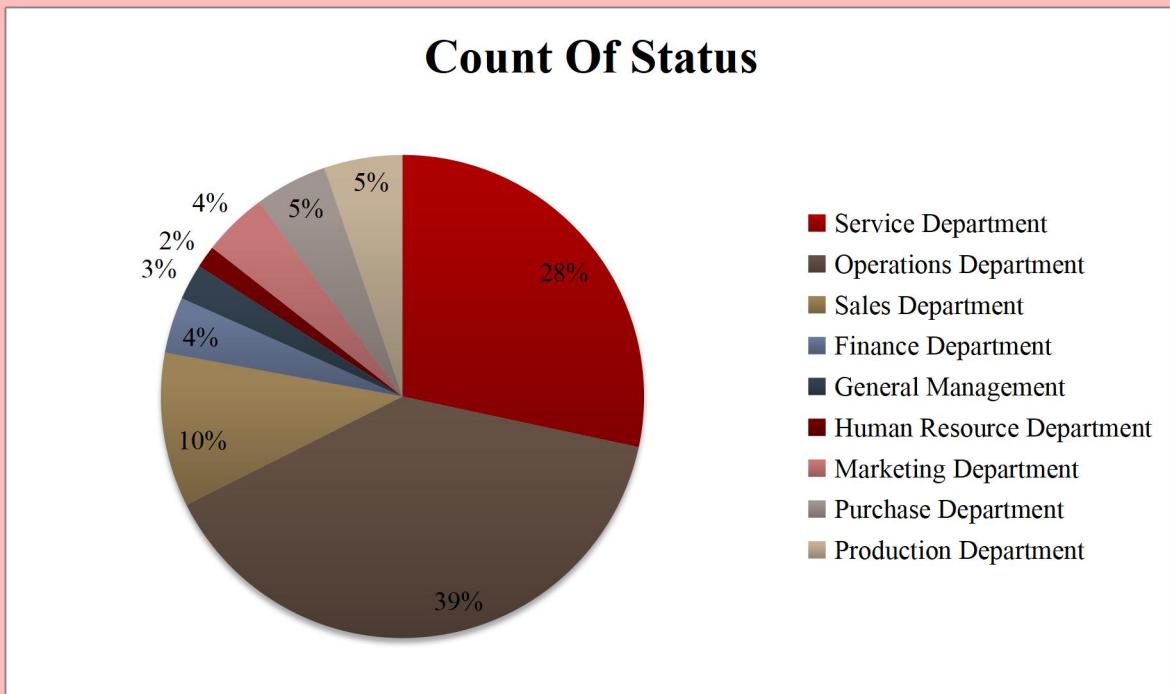
3. Charts and Plots: This is one of the most

important part of analysis to visualize the data.

Draw Pie Chart / Bar Graph ( or any other graph )

to show proportion of people working different department ?

Department	Count Of Status
Service Department	1332
Operations Department	1843
Sales Department	485
Finance Department	176
General Management	113
Human Resource Department	70
Marketing Department	202
Purchase Department	230
Production Department	246



From the above pie chart and bar plot I have analyzed  
that most of the people are working on the

Operations Department (i.e) people hired on this particular department and Least number of people hired for Human Resource Department.

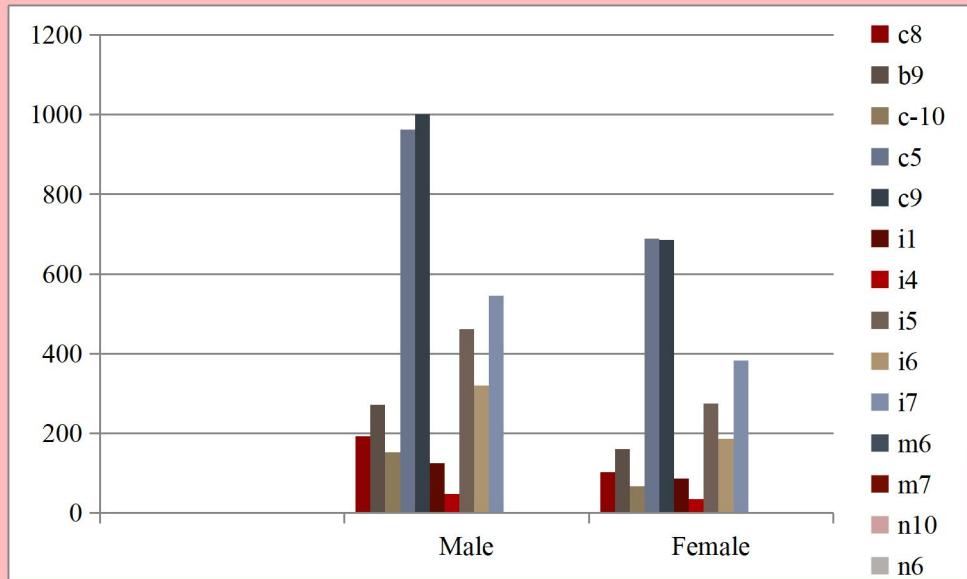
39% hired people that is ,(1843) people are working in the Operations Department and 1% of hired people that is, (70) only people working in the Human Resource Department.

#### **Findings 4**

**4. Charts:** Use different charts and graphs to perform the task representing the data.

**Represent different post tiers using chart/graph?**

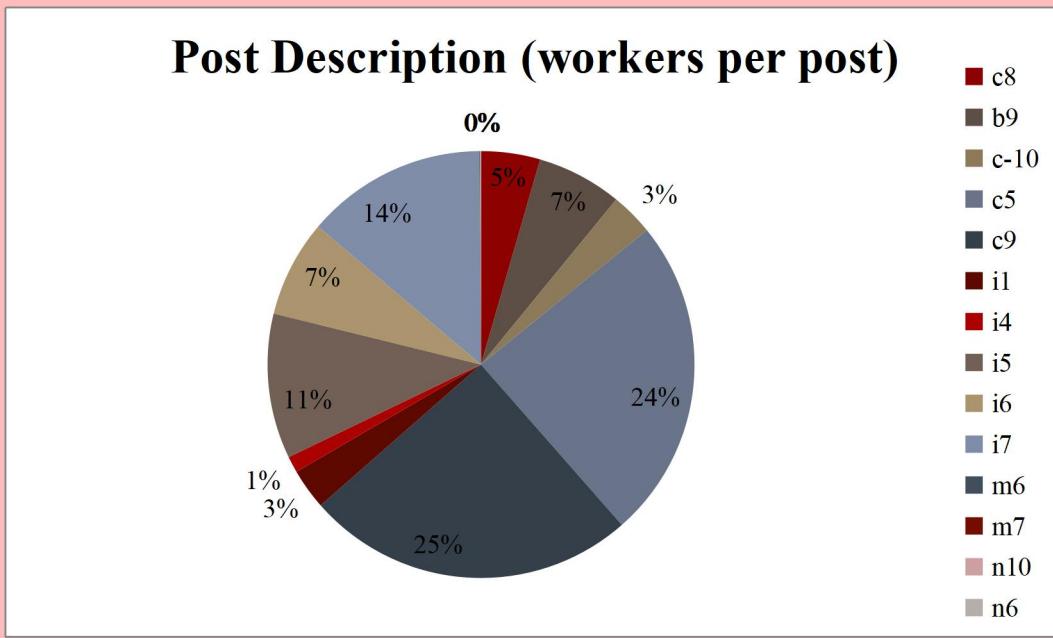
Event_name	c8	b9	c-10	c5	c9	i1	i4	i5	i6	i7	m6	m7	n10	n6	n9
Male	193	271	152	963	1001	125	48	461	320	546	2	1	1	0	0
Female	103	161	67	689	686	86	35	275	187	383	1	0	0	1	1



**From above column chart, Post Description (Gender wise distribution per post) is then male has the highest post tiers than the female employees.**

No. of People	post_name
320	c8
463	b9
232	c-10
1747	c5
1792	c9
222	i1
88	i4
787	i5
527	i6
982	i7
3	m6
1	m7
1	n10
1	n6
1	n9
1	-





From above pie chart , I have inferred a post description (workers per post) then there c9 post of 25% is the most workers working in c9 post.

### Findings 5

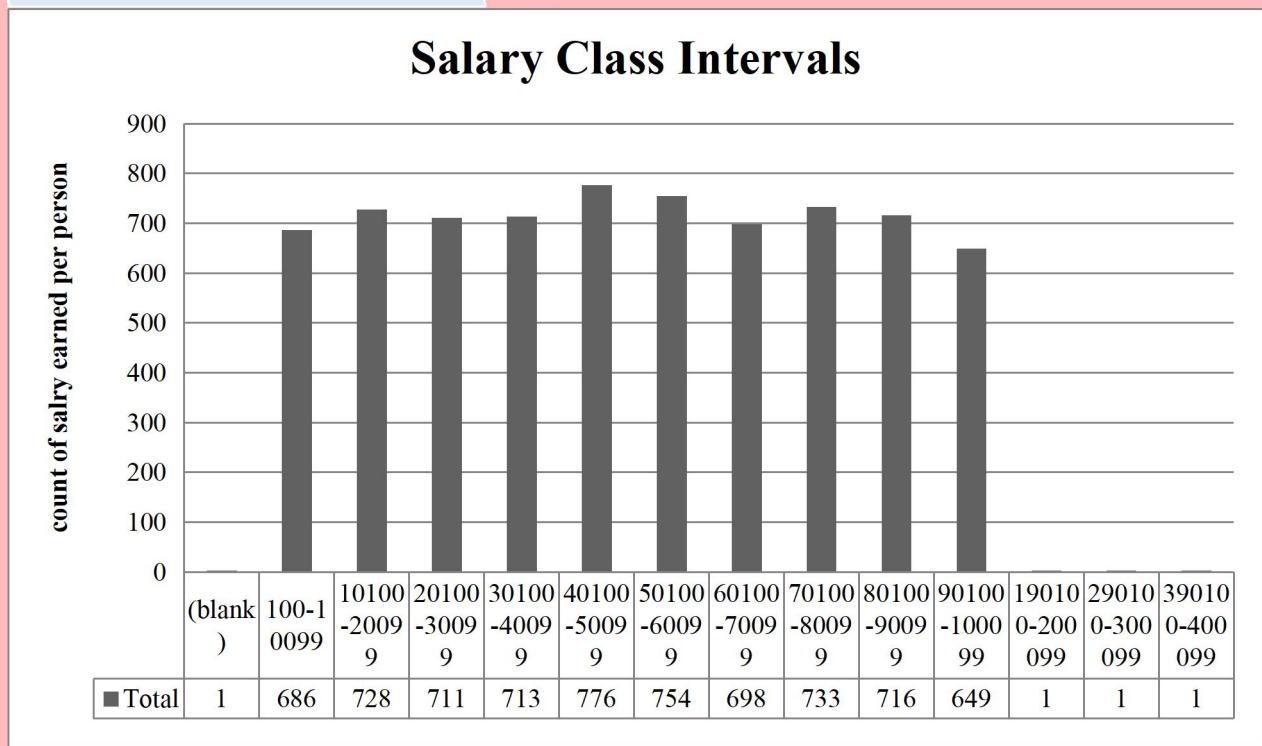
5. Class Intervals: The class interval is the difference between the upper class limit and the lower class limit.

Draw the class intervals for salary in the company ?

Row Labels	Count of Status
(blank)	1
100-10099	686
10100-20099	728
20100-30099	711
30100-40099	713
40100-50099	776
50100-60099	754
60100-70099	698
70100-80099	733
80100-90099	716
90100-100099	649
190100-200099	1
290100-300099	1
390100-400099	1
<b>Grand Total</b>	<b>7168</b>



### Salary Class Intervals

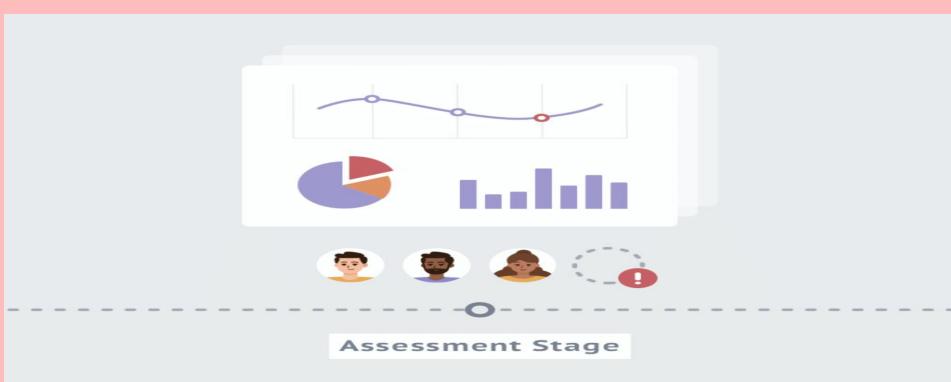


From the above class interval column chart represents that most of the people are on the intervals of **40100 – 50099** i.e, (**776**) are getting salary more than 40100.

## Conclusion

In the conclusion part, I would like to conclude that Hiring Process Analytics plays an important part for all the companies and firms to decide the job openings for the near future.

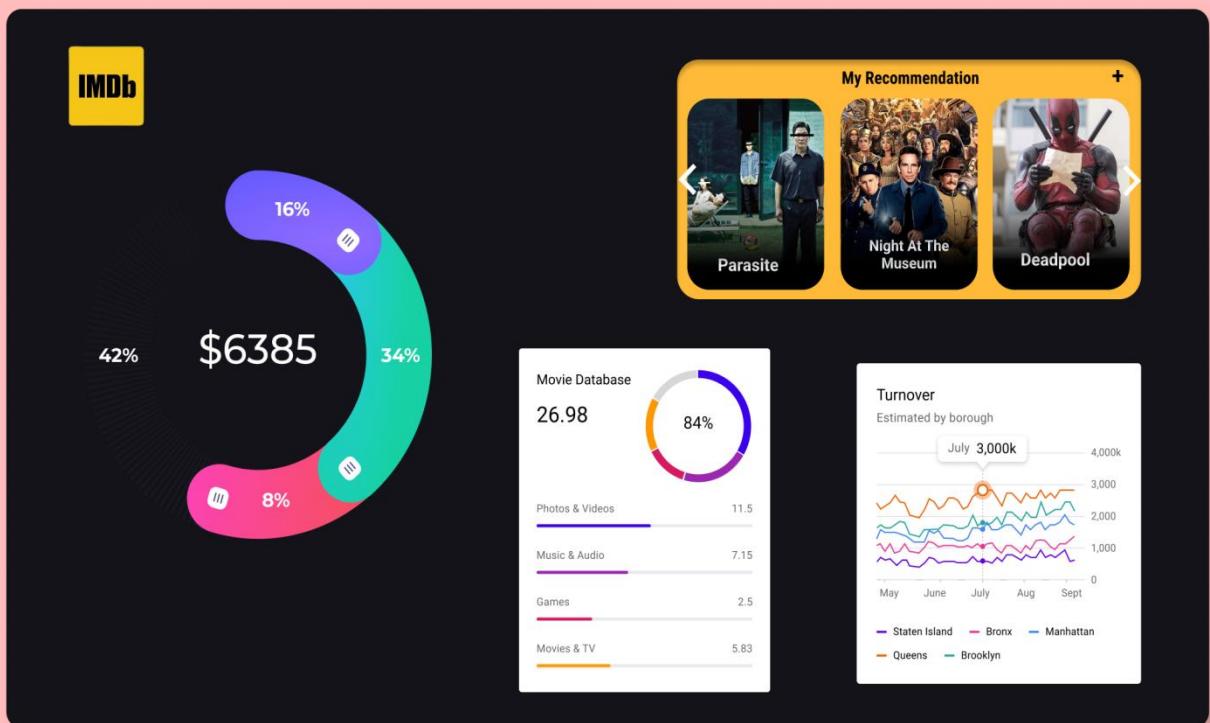
Hiring Process Analytics is done on monthly, quarterly or yearly basis as per the needs and policies of the companies. For any company the Operations Department has the highest number of workforce due to the workload on this department as this department acts as a central hub for all the executive tasks carried out. For any company there will some employees who have high salary packages compared to other employees, and this is due to the fact that they have some special skills and years of experience in their particular field of work. 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 employees.



# IMDb Movie Analysis

## Final Project-1

(Project – 5)



## **DESCRIPTION**

**Problem Statement:** The dataset provided is related to IMDB Movies. A potential problem to investigate could be: "What factors influence the success of a movie on IMDB?" Here, success can be defined by high IMDB ratings. The impact of this problem is significant for movie producers, directors, and investors who want to understand what makes a movie successful to make informed decisions in their future projects.

**Data Cleaning:** This step involves preprocessing the data to make it suitable for analysis. It includes handling missing values, removing duplicates, converting data types if necessary, and possibly feature engineering.

**Data Analysis:** Here, you'll explore the data to understand the relationships between different variables. You might look at the correlation between movie ratings and other factors like genre, director, budget, etc. You might also want to consider the year of release, the actors involved, and other relevant factors.

## The Problem

- ✓ Clean the data

(Dropping columns, removing null values, etc)

- ✓ Find the movies with the highest profit?

(observing the outliers using the appropriate chart type)

- ✓ Find IMDB Top 250

(Extract all the movies in the IMDb\_Top\_250 column which are not in the English language and store them in a new column named Top\_Foreign\_Lang\_Film.)

- ✓ Find the best directors

(Top 10 director by highest mean of IMDB score)

- ✓ Find popular genres

(Perform this by step while performing previous steps)

- ✓ Find the critic-favorite and audience-favorite actors

(Finding the mean of the num\_critic\_for\_reviews and num\_users\_for\_review and identify the actors which have the highest mean) and with decade of new data frame called df\_by\_decade.

## **Design**

- 1 . Firstly I made a copy of the raw data where I can perform the Analysis so that what ever changes I made it will not affect the original data.
2. Then dropping the columns which have no use for the analysis that we will be doing Columns like ‘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’ are the columns containing irrelevant data for the analysis tasks provided. So, these columns needs to be dropped.
3. After dropping the irrelevant columns now we need to remove the rows from the dataset having anyone of its column value as blank/NULL. Then we need to get rid off the duplicate values in the dataset which can be achieved by using the ‘Remove Duplicate Values/Cells’ available in the ‘Data’ tab

## Findings 1

1. **Cleaning the data:** This is one of the most important step to perform before moving forward with the analysis. Use your knowledge learned till now to do this.  
(Dropping columns, removing null values, etc.

### Your task: Clean the data

- ❖ First of all I have analyzed that with the provided dataset , there are 5043 rows and 28 columns and 1 header row are present there.
- ❖ After then, I looked a project description that related and required columns can only be used for the upcoming analysis so then I have been used a data cleaning strategy using MySQL workbench and excel functions to solve the data cleaning task.
- ❖ So I need only the columns and they are :

- **director\_name**
- **num\_critic\_for\_reviews**
- **gross**
- **genres**
- **actor\_1\_name**
- **movie\_title**
- **num\_voted\_users**
- **num\_user\_for\_reviews**

- **language**
- **budget**
- **title\_year**
- **imdb\_score**
- **movie\_facebook\_likes**
- ❖ After these are the only included then I left with a **13 columns** . Then I go with a null values and null rows , this was get rid by using the excel function (press F5 – go to special – select blanks – press ok) , it will show the blank rows and then ( data – remove duplicates) it will remove the duplicates and save this sheet . Then load this currently saved dataset in the MySQL workbench , queried for removing the extra null values precisely.
- ❖ And save the removed null values dataset and again load in the excel to see how many entries has been left .
- ❖ By performing these functions in the excel I found that the total number of rows is then 3789.
- ❖ Below this , I have attached with cleaned dataset as a picture of removing duplicates and with the blanks, required columns , and set the blank(language) column been default as ‘English’.
- ❖ This is my attached excel file.

[https://docs.google.com/spreadsheets/d/1Ha--2RrVp2mce6NS6hj5dGjqWq4jOpiw/edit?usp=drive\\_link](https://docs.google.com/spreadsheets/d/1Ha--2RrVp2mce6NS6hj5dGjqWq4jOpiw/edit?usp=drive_link)



## Findings 2

**2.Movies with highest profit:** Create a new column called profit which contains the difference of the two columns: gross and budget. Sort the column using the profit column as reference. Plot profit (y-axis) vs budget (x- axis) and observe the outliers using the appropriate chart type.

**Your task: Find the movies with the highest profit?**

- ❖ Now as the values in gross and budget columns are very large so I have used a excel to get rid of these large values is not readable that ([go to format cells – select custom – type #,,”million” – then select ok](#)) it will give us the readable values , i.e 274M.
- ❖ Next , we create a column called ‘profit’ as gross-budget difference in the MySQL workbench to get the desired results. The query can be as follows:

```
SELECT director_name,
       movie_title,
       ROUND(budget, 0) AS budget,
       ROUND(gross - budget, 0) AS profit
  FROM
    moviesp.cleanedm
 ORDER BY gross - budget DESC
 LIMIT 10;
```

With this query , I have found out the result of movie\_title are as follows:

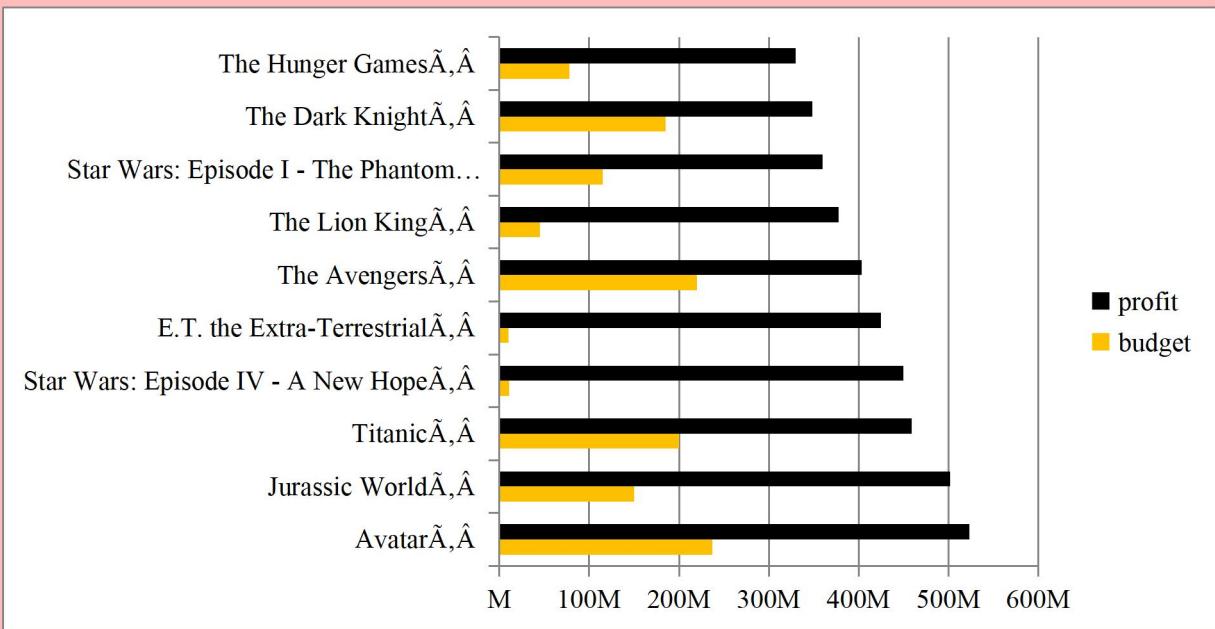
Avatar  
 Jurassic World  
 Titanic  
 Star Wars: Episode IV - A New Hope  
 E.T. the Extra-Terrestrial  
 The Avengers  
 The Lion King  
 Star Wars: Episode I - The Phantom Menace  
 The Dark Knight  
 The Hunger Games

director_name	movie_title	budget	profit
James Cameron	Avatar,Â	237M	524M
Colin Trevorrow	Jurassic World,Â	150M	502M
James Cameron	TitanicÂ,Â	200M	459M
George Lucas	Star Wars: Episode IV - A New HopeÂ,Â	11M	450M
Steven Spielberg	E.T. the Extra-TerrestrialÂ,Â	11M	424M
Joss Whedon	The AvengersÂ,Â	220M	403M
Roger Allers	The Lion KingÂ,Â	45M	378M
George Lucas	Star Wars: Episode I - The Phantom MenaceÂ,Â	115M	360M
Christopher Nolan	The Dark KnightÂ,Â	185M	348M
Gary Ross	The Hunger GamesÂ,Â	78M	330M

This the result of above query , that we have found that Movies with Highest profit . And the insights of this table are:

This is my attached excel file:

[https://docs.google.com/spreadsheets/d/1vU6T3pqE0nJI688PbSpdxeA8LBAjKr0R/edit?usp=drive\\_link](https://docs.google.com/spreadsheets/d/1vU6T3pqE0nJI688PbSpdxeA8LBAjKr0R/edit?usp=drive_link)



From the above bar chart represents that the highest profit of respected movies have been found to be analyzed and gave insights.

### Findings 3

**3.Top 250:** Create a new column IMDb\_Top\_250 and store the top 250 movies with the highest IMDb Rating (corresponding to the column: imdb\_score). Also make sure that for all of these movies, the num\_voted\_users is greater than 25,000. Also add a Rank column containing the values 1 to 250 indicating the ranks of the corresponding films.

Extract all the movies in the IMDb\_Top\_250 column which are not in the English language and store them in a new column named Top\_Foreign\_Lang\_Film. You can use your own imagination also!

### Your task: Find IMDB Top 250

a)

F

or the first part we should look for all the movies which has higher imdb scores user votes. So that I had solved this question by query with the user votes of higher than 30000 are as follows:

```
SELECT row_number() over(order by imdb_score DESC,
num_voted_users DESC) as ranking,
imdb_score, num_voted_users, movie_title AS IMDb_Top_250,
language
FROM moviesp.cleanedm
WHERE num_voted_users > 30000
LIMIT 250;
```

IMDB Top 250 Movies			
ranking	imdb_score	IMDb_Top_250	language
1	9.3	The Shawshank Redemption	English
2	9	The Dark Knight	English
3	9	The Godfather: Part II	English
4	8.9	The Lord of the Rings: The Return of the King	English
5	8.9	Schindler's List	English
6	8.8	Inception	English
7	8.8	Fight Club	English
8	8.8	Forrest Gump	English

9	8.8	The Lord of the Rings: The Fellowship of the Ring	English
10	8.8	Star Wars: Episode V - The Empire Strikes Back	English
11	8.7	The Matrix	English
12	8.7	The Lord of the Rings: The Two Towers	English
13	8.7	Star Wars: Episode IV - A New Hope	English
14	8.7	Goodfellas	English
15	8.6	Se7en	English
16	8.6	Interstellar	English
17	8.6	The Silence of the Lambs	English
18	8.6	Saving Private Ryan	English
19	8.6	American History X	English
20	8.6	Spirited Away	Japanese
21	8.5	The Dark Knight Rises	English
22	8.5	Gladiator	English
23	8.5	Django Unchained	English
24	8.5	The Departed	English
25	8.5	The Prestige	English
26	8.5	The Green Mile	English
27	8.5	Terminator 2: Judgment Day	English
28	8.5	Back to the Future	English
29	8.5	Raiders of the Lost Ark	English
30	8.5	The Lion King	English
31	8.5	The Pianist	English
32	8.5	Apocalypse Now	English
33	8.5	Psycho	English
34	8.4	American Beauty	English
35	8.4	Braveheart	English
36	8.4	WALL-E	English
37	8.4	Star Wars: Episode VI - Return of the Jedi	English
38	8.4	Amélie	French
39	8.4	Aliens	English
40	8.4	Once Upon a Time in America	English
41	8.4	Lawrence of Arabia	English
42	8.4	Das Boot	German
43	8.4	Baahubali: The Beginning	Telugu
44	8.3	Batman Begins	English
45	8.3	Inglourious Basterds	English
46	8.3	Eternal Sunshine of the Spotless Mind	English
47	8.3	Up	English
48	8.3	Toy Story	English
49	8.3	Good Will Hunting	English
50	8.3	Snatch	English
51	8.3	Toy Story 3	English
52	8.3	Scarface	English
53	8.3	Indiana Jones and the Last Crusade	English
54	8.3	2001: A Space Odyssey	English

55	8.3	L.A. ConfidentialÃ¢ÂÂ	English
56	8.3	Inside OutÃ¢ÂÂ	English
57	8.3	UnforgivenÃ¢ÂÂ	English
58	8.3	AmadeusÃ¢ÂÂ	English
59	8.3	DownfallÃ¢ÂÂ	German
60	8.3	Raging BullÃ¢ÂÂ	English
61	8.3	RoomÃ¢ÂÂ	English
62	8.3	MetropolisÃ¢ÂÂ	German
63	8.2	V for VendettaÃ¢ÂÂ	English
64	8.2	The Wolf of Wall StreetÃ¢ÂÂ	English
65	8.2	Finding NemoÃ¢ÂÂ	English
66	8.2	A Beautiful MindÃ¢ÂÂ	English
67	8.2	Die HardÃ¢ÂÂ	English
68	8.2	Gran TorinoÃ¢ÂÂ	English
69	8.2	The Big LebowskiÃ¢ÂÂ	English
70	8.2	How to Train Your DragonÃ¢ÂÂ	English
71	8.2	Pan's LabyrinthÃ¢ÂÂ	Spanish
72	8.2	Blade RunnerÃ¢ÂÂ	English
73	8.2	Into the WildÃ¢ÂÂ	English
74	8.2	CasinoÃ¢ÂÂ	English
75	8.2	WarriorÃ¢ÂÂ	English
76	8.2	Captain America: Civil WarÃ¢ÂÂ	English
77	8.2	The ThingÃ¢ÂÂ	English
78	8.2	Howl's Moving CastleÃ¢ÂÂ	Japanese
79	8.1	The AvengersÃ¢ÂÂ	English
80	8.1	Pirates of the Caribbean: The Curse of the Black PearlÃ¢ÂÂ	English
81	8.1	Shutter IslandÃ¢ÂÂ	English
82	8.1	Kill Bill: Vol. 1Ã¢ÂÂ	English
83	8.1	The Sixth SenseÃ¢ÂÂ	English
84	8.1	Guardians of the GalaxyÃ¢ÂÂ	English
85	8.1	The Truman ShowÃ¢ÂÂ	English
86	8.1	Sin CityÃ¢ÂÂ	English
87	8.1	Jurassic ParkÃ¢ÂÂ	English
88	8.1	No Country for Old MenÃ¢ÂÂ	English
89	8.1	Monsters, Inc.Ã¢ÂÂ	English
90	8.1	Gone GirlÃ¢ÂÂ	English
91	8.1	Mad Max: Fury RoadÃ¢ÂÂ	English
92	8.1	The Bourne UltimatumÃ¢ÂÂ	English
93	8.1	Million Dollar BabyÃ¢ÂÂ	English
94	8.1	DeadpoolÃ¢ÂÂ	English
95	8.1	The Grand Budapest HotelÃ¢ÂÂ	English
96	8.1	The MartianÃ¢ÂÂ	English
97	8.1	The Imitation GameÃ¢ÂÂ	English
98	8.1	12 Years a SlaveÃ¢ÂÂ	English
99	8.1	Groundhog DayÃ¢ÂÂ	English
100	8.1	The RevenantÃ¢ÂÂ	English

101	8.1	PrisonersÃ‚Â	English
102	8.1	There Will Be BloodÃ‚Â	English
103	8.1	The HelpÃ‚Â	English
104	8.1	RushÃ‚Â	English
105	8.1	The Princess BrideÃ‚Â	English
106	8.1	Hotel RwandaÃ‚Â	English
107	8.1	SpotlightÃ‚Â	English
108	8.1	The Sea InsideÃ‚Â	Spanish
109	8.1	Tae Guk Gi: The Brotherhood of WarÃ‚Â	Korean
110	8	Slumdog MillionaireÃ‚Â	English
111	8	Black SwanÃ‚Â	English
112	8	District 9Ã‚Â	English
113	8	Catch Me If You CanÃ‚Â	English
114	8	X-Men: Days of Future PastÃ‚Â	English
115	8	Kill Bill: Vol. 2Ã‚Â	English
116	8	Star TrekÃ‚Â	English
117	8	The King's SpeechÃ‚Â	English
118	8	The IncrediblesÃ‚Â	English
119	8	RatatouilleÃ‚Â	English
120	8	Casino RoyaleÃ‚Â	English
121	8	Casino RoyaleÃ‚Â	English
122	8	Life of PiÃ‚Â	English
123	8	JawsÃ‚Â	English
124	8	Blood DiamondÃ‚Â	English
125	8	Rain ManÃ‚Â	English
126	8	HerÃ‚Â	English
127	8	The Perks of Being a WallflowerÃ‚Â	English
128	8	Big FishÃ‚Â	English
129	8	Mystic RiverÃ‚Â	English
130	8	The Pursuit of HappynessÃ‚Â	English
131	8	In BrugesÃ‚Â	English
132	8	The ExorcistÃ‚Â	English
133	8	Dead Poets SocietyÃ‚Â	English
134	8	AladdinÃ‚Â	English
135	8	SerenityÃ‚Â	English
136	8	MagnoliaÃ‚Â	English
137	8	Mulholland DriveÃ‚Â	English
138	8	The ArtistÃ‚Â	English
139	8	Dances with WolvesÃ‚Â	English
140	8	True RomanceÃ‚Â	English
141	8	BrazilÃ‚Â	English
142	8	Cinderella ManÃ‚Â	English
143	8	The Iron GiantÃ‚Â	English
144	8	JFKÃ‚Â	English
145	8	Dancer in the DarkÃ‚Â	English
146	8	Doctor ZhivagoÃ‚Â	English

147	7.9	Avatar	English
148	7.9	Iron Man	English
149	7.9	The Hobbit: An Unexpected Journey	English
150	7.9	Taken	English
151	7.9	The Hobbit: The Desolation of Smaug	English
152	7.9	Shrek	English
153	7.9	Edge of Tomorrow	English
154	7.9	The Bourne Identity	English
155	7.9	The Notebook	English
156	7.9	Toy Story 2	English
157	7.9	Children of Men	English
158	7.9	Edward Scissorhands	English
159	7.9	Hot Fuzz	English
160	7.9	Captain Phillips	English
161	7.9	E.T. the Extra-Terrestrial	English
162	7.9	Big Hero 6	English
163	7.9	The Fighter	English
164	7.9	The Hateful Eight	English
165	7.9	How to Train Your Dragon 2	English
166	7.9	The Untouchables	English
167	7.9	Crouching Tiger, Hidden Dragon	Mandarin
168	7.9	Almost Famous	English
169	7.9	Boogie Nights	English
170	7.9	Walk the Line	English
171	7.9	Halloween	English
172	7.9	Hero	Mandarin
173	7.9	The Blues Brothers	English
174	7.9	Ed Wood	English
175	7.9	The Insider	English
176	7.9	Letters from Iwo Jima	Japanese
177	7.9	Straight Outta Compton	English
178	7.9	Glory	English
179	7.9	Glory	English
180	7.9	My Fair Lady	English
181	7.9	The Remains of the Day	English
182	7.9	The Right Stuff	English
183	7.9	The World's Fastest Indian	English
184	7.8	The Hangover	English
185	7.8	Gravity	English
186	7.8	Silver Linings Playbook	English
187	7.8	Skyfall	English
188	7.8	X-Men: First Class	English
189	7.8	Captain America: The Winter Soldier	English
190	7.8	The Curious Case of Benjamin Button	English
191	7.8	Drive	English
192	7.8	Ocean's Eleven	English

193	7.8	Star Trek Into Darkness	English
194	7.8	Birdman or (The Unexpected Virtue of Ignorance)	English
195	7.8	Harry Potter and the Prisoner of Azkaban	English
196	7.8	The Bourne Supremacy	English
197	7.8	Back to the Future Part II	English
198	7.8	The Girl with the Dragon Tattoo	English
199	7.8	American Gangster	English
200	7.8	Tangled	English
201	7.8	Predator	English
202	7.8	Wreck-It Ralph	English
203	7.8	Lucky Number Slevin	English
204	7.8	The Game	English
205	7.8	Being John Malkovich	English
206	7.8	The Fault in Our Stars	English
207	7.8	The Lego Movie	English
208	7.8	3:10 to Yuma	English
209	7.8	Moonrise Kingdom	English
210	7.8	Apocalypto	Maya
211	7.8	Donnie Brasco	English
212	7.8	O Brother, Where Art Thou?	English
213	7.8	Gattaca	English
214	7.8	The Fugitive	English
215	7.8	About Time	English
216	7.8	Office Space	English
217	7.8	Changeling	English
218	7.8	Pride & Prejudice	English
219	7.8	Atonement	English
220	7.8	The Big Short	English
221	7.8	Finding Neverland	English
222	7.8	What's Eating Gilbert Grape	English
223	7.8	South Park: Bigger Longer & Uncut	English
224	7.8	Remember the Titans	English
225	7.8	Fantastic Mr. Fox	English
226	7.8	The Boy in the Striped Pajamas	English
227	7.8	The Last of the Mohicans	English
228	7.8	The Jungle Book	English
229	7.8	The Jungle Book	English
230	7.8	Kung Fu Hustle	Cantonese
231	7.8	Tombstone	English
232	7.8	Nebraska	English
233	7.8	Glengarry Glen Ross	English
234	7.8	The Last Emperor	English
235	7.8	The Best Offer	English
236	7.8	The Conjuring 2	English
237	7.8	The Color Purple	English
238	7.8	Black Book	Dutch

239	7.8	The White Ribbon	German
240	7.8	Hamlet	English
241	7.8	The Little Prince	English
242	7.8	The Verdict	English
243	7.7	Titanic	English
244	7.7	300	English
245	7.7	The Social Network	English
246	7.7	Argo	English
247	7.7	Kick-Ass	English
248	7.7	Minority Report	English
249	7.7	Cast Away	English
250	7.7	Watchmen	English

From the above query , I got this **IMDb\_Top\_250** , I had attached as proof of this.

This is my attached excel file:

[https://docs.google.com/spreadsheets/d/16NnCBQPx46wDCiuNatK76Pn8Ra47QAzb/edit?usp=drive\\_link](https://docs.google.com/spreadsheets/d/16NnCBQPx46wDCiuNatK76Pn8Ra47QAzb/edit?usp=drive_link)

b)For the next part as , out of these top250 movies with a highest user votes , we search for all the **Top\_Foreign\_Lang\_Film** using the below query :

```
SELECT row_number() over(order by imdb_score DESC,
num_voted_users DESC) as ranking,
imdb_score,
num_voted_users,
movie_title AS Top_non_english,
language
FROM
moviesp.cleanedm
```

**WHERE**

**num\_voted\_users > 30000 and language <> 'English'**

**LIMIT 36;**

ranking	imdb_score	num_voted_users	Top_non_english	language
1	8.6	417971	Spirited Away	Japanese
2	8.4	534262	Amélie	French
3	8.4	168203	Das Boot	German
4	8.4	62756	Baahubali: The Beginning	Telugu
5	8.3	248354	Downfall	German
6	8.3	111841	Metropolis	German
7	8.2	467234	Pan's Labyrinth	Spanish
8	8.2	214091	Howl's Moving Castle	Japanese
9	8.1	64556	The Sea Inside	Spanish
10	8.1	31943	Tae Guk Gi: The Brotherhood of War	Korean
11	7.9	217740	Crouching Tiger, Hidden Dragon	Mandarin
12	7.9	149414	Hero	Mandarin
13	7.9	132149	Letters from Iwo Jima	Japanese
14	7.8	236000	Apocalypto	Maya
15	7.8	99353	Kung Fu Hustle	Cantonese
16	7.8	59507	Black Book	Dutch
17	7.8	52958	The White Ribbon	German
18	7.7	85589	Ponyo	Japanese
19	7.7	62607	A Very Long Engagement	French
20	7.7	55516	The Great Beauty	Italian
21	7.6	92295	House of Flying Daggers	Mandarin
22	7.6	68119	The Kite Runner	Dari
23	7.6	38690	The Flowers of War	Mandarin
24	7.5	77656	The Painted Veil	Mandarin
25	7.5	41496	2046	Cantonese
26	7.4	36894	Red Cliff	Mandarin
27	7.3	63084	Paris, je t'aime	French
28	7.3	37635	Mongol: The Rise of Genghis Khan	Mongolian
29	7.2	55928	District B13	French
30	7.1	179235	The Passion of the Christ	Aramaic
31	7	36455	Curse of the Golden Flower	Mandarin
32	6.7	32003	Coco Before Chanel	French
33	6.4	86152	The Interpreter	Aboriginal
34	5.9	71574	The Legend of Zorro	Spanish
35	4.3	31414	In the Land of Blood and Honey	Bosnian

- ❖ From the above query , I have found **35 movies** list with respect to the **num\_voted\_users** for non English movies with the top 250 lists.
- ❖ Then with this results , have been queried with respect to the number of voters for the most favorable movies with ratings of **imdb\_score**, that is **non\_english** that is foreign language movies that have been listed using the query. So in this question Top 250 and **Top\_Foreign\_Lang\_Film(35 movies)** , MySQL workbench had been used to get the desired results and with the results , I had transferred the results to excel to make table to be readable.
- ❖ This is my file attached here:  
[https://drive.google.com/file/d/1sDk5gBzVoLVdENnGPDrv77P4AksNeHbW/view?usp=drive\\_link](https://drive.google.com/file/d/1sDk5gBzVoLVdENnGPDrv77P4AksNeHbW/view?usp=drive_link)

## **Findings 4**

- 1. Best Directors:** Group the column using the **director\_name** column.

Find out the top 10 directors for whom the mean of **imdb\_score** is the highest and store them in a new column **top10director**. In case of a tie in IMDb score between two directors, sort them alphabetically.

## **Your task: Find the best directors**

For this question ,want to find out the top 10 directors whom with the imdb\_score of highest one . This will be found out by using the below query:

```
SELECT
    director_name AS Top_10_Director,
    ROUND(AVG(imdb_score), 2) AS avg_rating
FROM
    moviesp.cleanedm
GROUP BY director_name
ORDER BY avg_rating DESC , director_name
LIMIT 10;
```

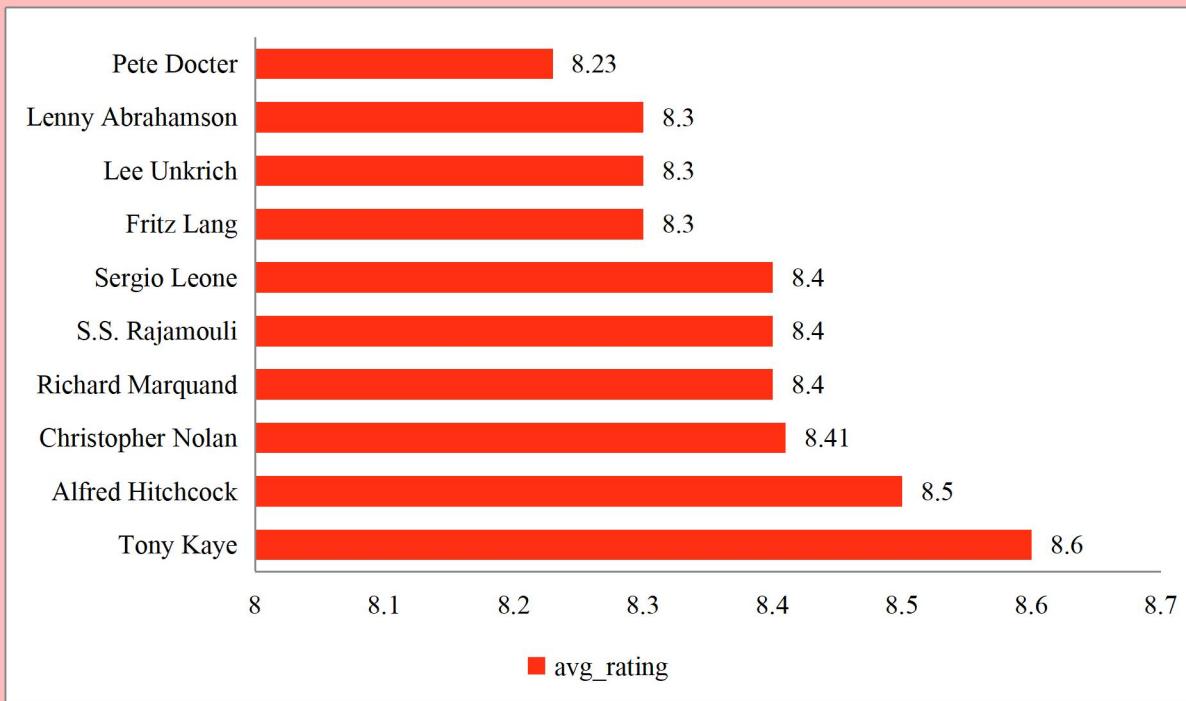
With this query , I have got the results.....

Top 10 Directors	
Top_10_Director	avg_rating
Tony Kaye	8.6
Alfred Hitchcock	8.5
Christopher Nolan	8.41
Richard Marquand	8.4
S.S. Rajamouli	8.4
Sergio Leone	8.4
Fritz Lang	8.3
Lee Unkrich	8.3
Lenny Abrahamson	8.3
Pete Docter	8.23

From the above results , I have inferred that the top 10 directors have been listed by using the sql query results with the below insights.

This is my excel file attached:

[https://docs.google.com/spreadsheets/d/1kDmzbNlehAevghQAK8UPAuYi1rA3jIXj/edit?usp=drive\\_link](https://docs.google.com/spreadsheets/d/1kDmzbNlehAevghQAK8UPAuYi1rA3jIXj/edit?usp=drive_link)



## Findings 5

2. **Popular Genres:** Perform this step using the knowledge gained while performing previous steps.

**Your task: Find popular genres**

- ❖ For this question, we first split the genres of each movie into individual genres into 8 separate levels of genres. Then, we make a list of all genres by uniting all

levels of genres and then count the popularity of each genre.

- ❖ I have been using sql query to solve this description of separating 8 levels of genres with list of all genres the with mean of gross and the number of movies.

**With d8 as**

```
(SELECT *, (CASE WHEN SUBSTRING_INDEX(SUBSTRING_INDEX(genres, '|', 8), '|', -1) = genre_1 THEN NULL
WHEN SUBSTRING_INDEX(SUBSTRING_INDEX(genres, '|', 8), '|', -1) = genre_2 THEN NULL
WHEN SUBSTRING_INDEX(SUBSTRING_INDEX(genres, '|', 8), '|', -1) = genre_3 THEN NULL
WHEN SUBSTRING_INDEX(SUBSTRING_INDEX(genres, '|', 8), '|', -1) = genre_4 THEN NULL
WHEN SUBSTRING_INDEX(SUBSTRING_INDEX(genres, '|', 8), '|', -1) = genre_5 THEN NULL
WHEN SUBSTRING_INDEX(SUBSTRING_INDEX(genres, '|', 8), '|', -1) = genre_6 THEN NULL
WHEN SUBSTRING_INDEX(SUBSTRING_INDEX(genres, '|', 8), '|', -1) = genre_7 THEN NULL
ELSE SUBSTRING_INDEX(SUBSTRING_INDEX(genres, '|', 8), '|', -1) END) AS genre_8
FROM
(SELECT *, (CASE WHEN SUBSTRING_INDEX(SUBSTRING_INDEX(genres, '|', 7), '|', -1) = genre_1 THEN NULL
WHEN SUBSTRING_INDEX(SUBSTRING_INDEX(genres, '|', 7), '|', -1) = genre_2 THEN NULL
WHEN SUBSTRING_INDEX(SUBSTRING_INDEX(genres, '|', 7), '|', -1) = genre_3 THEN NULL
WHEN SUBSTRING_INDEX(SUBSTRING_INDEX(genres, '|', 7), '|', -1) = genre_4 THEN NULL
WHEN SUBSTRING_INDEX(SUBSTRING_INDEX(genres, '|', 7), '|', -1) = genre_5 THEN NULL
WHEN SUBSTRING_INDEX(SUBSTRING_INDEX(genres, '|', 7), '|', -1) = genre_6 THEN NULL
ELSE SUBSTRING_INDEX(SUBSTRING_INDEX(genres, '|', 7), '|', -1) END) AS genre_7
FROM
(SELECT *, (CASE WHEN SUBSTRING_INDEX(SUBSTRING_INDEX(genres, '|', 6), '|', -1) = genre_1 THEN NULL
WHEN SUBSTRING_INDEX(SUBSTRING_INDEX(genres, '|', 6), '|', -1) = genre_2 THEN NULL
WHEN SUBSTRING_INDEX(SUBSTRING_INDEX(genres, '|', 6), '|', -1) = genre_3 THEN NULL
WHEN SUBSTRING_INDEX(SUBSTRING_INDEX(genres, '|', 6), '|', -1) = genre_4 THEN NULL
WHEN SUBSTRING_INDEX(SUBSTRING_INDEX(genres, '|', 6), '|', -1) = genre_5 THEN NULL
ELSE SUBSTRING_INDEX(SUBSTRING_INDEX(genres, '|', 6), '|', -1) END) AS genre_6
FROM
(SELECT *, (CASE WHEN SUBSTRING_INDEX(SUBSTRING_INDEX(genres, '|', 5), '|', -1) = genre_1 THEN NULL
WHEN SUBSTRING_INDEX(SUBSTRING_INDEX(genres, '|', 5), '|', -1) = genre_2 THEN NULL
WHEN SUBSTRING_INDEX(SUBSTRING_INDEX(genres, '|', 5), '|', -1) = genre_3 THEN NULL
WHEN SUBSTRING_INDEX(SUBSTRING_INDEX(genres, '|', 5), '|', -1) = genre_4 THEN NULL
ELSE SUBSTRING_INDEX(SUBSTRING_INDEX(genres, '|', 5), '|', -1) END) AS genre_5
FROM(SELECT *, (CASE WHEN SUBSTRING_INDEX(SUBSTRING_INDEX(genres, '|', 4), '|', -1) = genre_1 THEN NULL
WHEN SUBSTRING_INDEX(SUBSTRING_INDEX(genres, '|', 4), '|', -1) = genre_2 THEN NULL
WHEN SUBSTRING_INDEX(SUBSTRING_INDEX(genres, '|', 4), '|', -1) = genre_3 THEN NULL
ELSE SUBSTRING_INDEX(SUBSTRING_INDEX(genres, '|', 4), '|', -1) END) AS genre_4
FROM
(SELECT *,(CASE WHEN SUBSTRING_INDEX(SUBSTRING_INDEX(genres, '|', 3), '|', -1) = genre_1 THEN NULL
WHEN SUBSTRING_INDEX(SUBSTRING_INDEX(genres, '|', 3), '|', -1) = genre_2 THEN NULL
ELSE SUBSTRING_INDEX(SUBSTRING_INDEX(genres, '|', 3), '|', -1) END) AS genre_3
```

```

FROM
(SELECT *, (CASE WHEN SUBSTRING_INDEX(SUBSTRING_INDEX(genres, '|', 2), '|', -1) =
genre_1 THEN NULL
ELSE SUBSTRING_INDEX(SUBSTRING_INDEX(genres, '|', 2), '|', -1) END) AS genre_2
FROM
(SELECT
SUBSTRING(genres, 2) as genres,
SUBSTRING_INDEX(SUBSTRING_INDEX(SUBSTRING_INDEX(SUBSTRING(genres, 2), '|', 1), '|', -1) AS genre_1
FROM moviesp.cleanedm) d1 d2 d3 d4 d5 d6 d7),
list as
(SELECT DISTINCT genre_1 AS genre_list
FROM
(
(SELECT genre_1 FROM d8) UNION ALL
(SELECT genre_2 FROM d8) UNION ALL
(SELECT genre_3 FROM d8) UNION ALL
(SELECT genre_4 FROM d8) UNION ALL
(SELECT genre_5 FROM d8) UNION ALL
(SELECT genre_6 FROM d8) UNION ALL
(SELECT genre_7 FROM d8) UNION ALL
(SELECT genre_8 FROM d8)) g12
WHERE genre_1 IS NOT NULL ORDER BY genre_1)
SELECT genre_list,round(AVG(DISTINCT CASE WHEN INSTR(genres, genre_list) THEN
gross ELSE NULL END), 2) AS avg_grossing,
COUNT(DISTINCT CASE WHEN INSTR(genres, genre_list) THEN genres ELSE NULL END)
AS no_of_movies
FROM cleanedm JOIN list
GROUP BY genre_list
ORDER BY avg_grossing DESC;

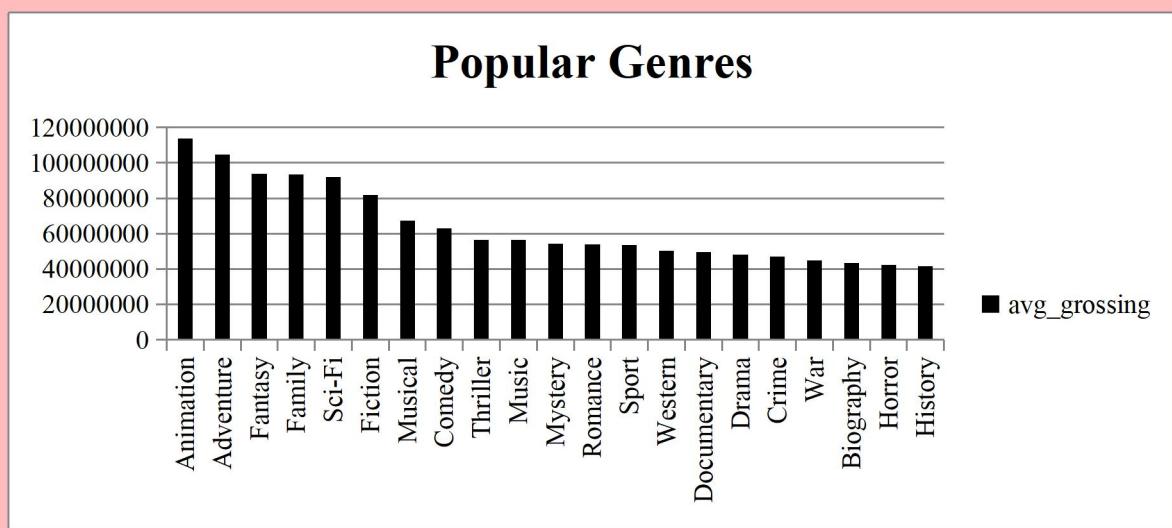
SELECT actor_1_name, COUNT(movie_title) AS no_of_movies,
ROUND(AVG(num_user_for_reviews), 2) AS user_reviews,
ROUND(AVG(num_critic_for_reviews), 2) AS critic_reviews
FROM
(
(SELECT actor_1_name, movie_title,num_user_for_reviews, num_critic_for_reviews
FROM moviesp.cleanedm
WHERE actor_1_name = ' CCH Pounder')
UNION ALL
(SELECT actor_1_name, movie_title,num_user_for_reviews, num_critic_for_reviews
FROM moviesp.cleanedm
WHERE actor_1_name = ' J.K. Simmons')
UNION ALL
(SELECT actor_1_name, movie_title, num_user_for_reviews, num_critic_for_reviews
FROM moviesp.cleanedm
WHERE actor_1_name = ' Jennifer Lawrence')) c
GROUP BY actor_1_name
ORDER BY user_reviews DESC , critic_reviews DESC;

```

Popular Genres		
genre_list	avg_grossing	no_of_movies
Animation	113547233.2	84
Adventure	104744835.7	238
Fantasy	93680359.4	170
Family	93247050.7	178

Sci-Fi	91891509.05	134
Fiction	81796457.68	235
Musical	67450846.54	53
Comedy	62876377.93	263
Thriller	56544554.87	179
Music	56385988.68	101
Mystery	54357927.28	97
Romance	53797743.97	171
Sport	53731204.96	39
Western	50326576.6	33
Documentary	49478940.43	4
Drama	47978352.15	310
Crime	46975720.56	121
War	44756370.57	60
Biography	43329712.5	63
Horror	42187374.12	70
History	41782106.2	51

From the above query results , There are **23 Genres** got to inferred that the splitted genres with the mean of gross and number of movies match lists. And got the analysis of insights of column chart are as follows:



This is my attached excel file:

[https://docs.google.com/spreadsheets/d/13gHz7H3oDULFNXcZJIRAB2g6VOb-UgVZ/edit?usp=drive\\_link](https://docs.google.com/spreadsheets/d/13gHz7H3oDULFNXcZJIRAB2g6VOb-UgVZ/edit?usp=drive_link)

## Findings 6

3. **Charts:** Create three new columns namely, **Meryl\_Streep, Leo\_Caprio, and Brad\_Pitt** which contain the movies in which the actors: 'Meryl Streep', 'Leonardo DiCaprio', and 'Brad Pitt' are the lead actors. Use only the **actor\_1\_name** column for extraction. Also, make sure that you use the names 'Meryl Streep', 'Leonardo DiCaprio', and 'Brad Pitt' for the said extraction.

Append the rows of all these columns and store them in a new column named **Combined**.

Group the combined column using the **actor\_1\_name** column.

Find the mean of the **num\_critic\_for\_reviews** and **num\_users\_for\_review** and identify the actors which have the highest mean.

Observe the change in number of voted users over decades using a bar chart. Create a column called **decade** which represents the decade to which every movie belongs to. For example, the **title\_year** year 1923, 1925 should be stored as

1920s. Sort the column based on the column decade, group it by decade and find the sum of users voted in each decade. Store this in a new data frame called `df_by_decade`.

### **Your task: Find the critic-favorite and audience-favorite actors**

a) For this critic-favourite , the assigned favourite movies has been extracted which contain the movies in which the actors: 'Meryl Streep', 'Leonardo DiCaprio', and 'Brad Pitt' are the lead actors.

❖ So I have queried for extracting the provided favourite actors with the number of movies , user reviews (audience-favourite) and critic-reviews. The query for finding this as follows:

```
SELECT actor_1_name, COUNT(movie_title) AS no_of_movies,  
ROUND(AVG(num_user_for_reviews), 2) AS user_reviews,  
ROUND(AVG(num_critic_for_reviews), 2) AS critic_reviews  
FROM  
(  
(SELECT actor_1_name, movie_title,num_user_for_reviews,  
num_critic_for_reviews  
FROM moviesp.cleanedm  
WHERE actor_1_name = 'Meryl Streep')  
UNION ALL
```

```

(SELECT actor_1_name, movie_title, num_user_for_reviews,
num_critic_for_reviews
FROM moviesp.cleanedm
WHERE actor_1_name = 'Leonardo DiCaprio')
UNION ALL
(SELECT actor_1_name, movie_title, num_user_for_reviews,
num_critic_for_reviews
FROM moviesp.cleanedm
WHERE actor_1_name = 'Brad Pitt')) c
GROUP BY actor_1_name
ORDER BY user_reviews DESC , critic_reviews DESC;

```

Favourite actors			
actor_1_name	num_of_movies	user_reviews	critic_reviews
Leonardo DiCaprio	20	922.55	322.2
Brad Pitt	17	742.35	245
Meryl Streep	11	297.18	181.45

- ❖ From the above query , I got this results which inferred of extracting from the actor name 1 with the help of user reviews and critic reviews.
  - ❖ And this query helped me to find out that Leonardo DiCaprio is the most favorite actor for both the audience and the critics.
  - ❖ This is the file attached here:
- [https://docs.google.com/spreadsheets/d/1046zg\\_F9vOST467Ewu6HOAwdWET-C-Bn/edit?usp=drive\\_link](https://docs.google.com/spreadsheets/d/1046zg_F9vOST467Ewu6HOAwdWET-C-Bn/edit?usp=drive_link)

b) Sort the column based on the column decade, group it by decade and find the sum of users voted in each decade. Store this in a new data frame called **df\_by\_decade**.

❖ For this decades per number user votes , we will bw finding out extracting the years from MySQL query workbench are as follows:

```
SELECT
    CONCAT(CONVERT( FLOOR(title_year / 10) * 10 ,
CHAR),
    's') AS decade,
    SUM(num_voted_users) AS total_votes
FROM
    moviesp.cleanedm
GROUP BY decade
ORDER BY decade;
```

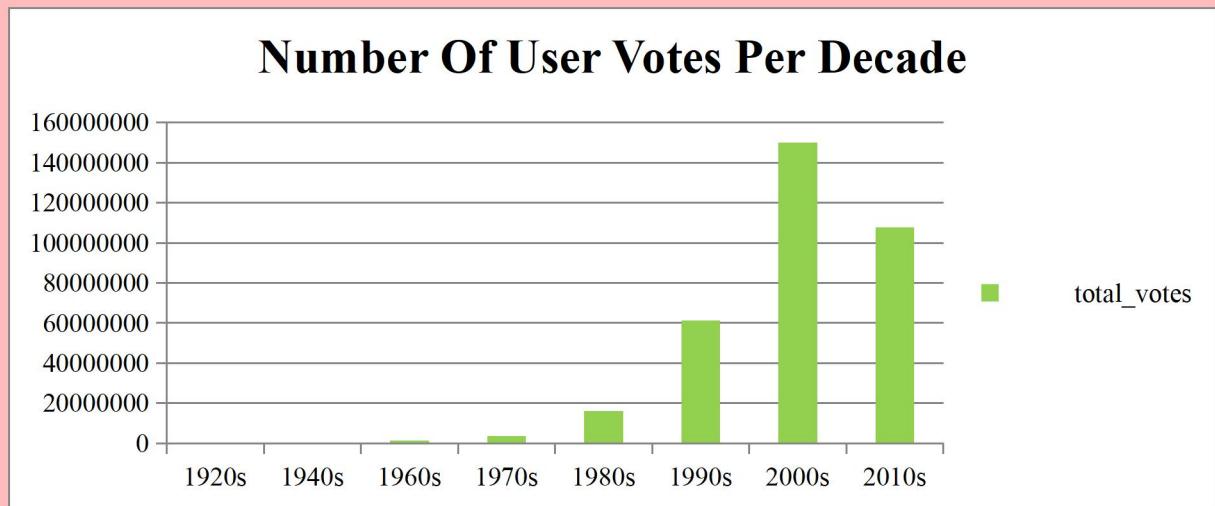
This is excel attached here:

[https://docs.google.com/spreadsheets/d/1B3LzOFwL59h8d2jBRb6YJ73Y3tmHTAiR/edit?usp=drive\\_link](https://docs.google.com/spreadsheets/d/1B3LzOFwL59h8d2jBRb6YJ73Y3tmHTAiR/edit?usp=drive_link)

<b>Number Of User Votes Per Decade</b>	
<b>decade</b>	<b>total_votes</b>
1920s	111841

1940s	90360
1960s	1218680
1970s	3679621
1980s	16136813
1990s	61213572
2000s	150036196
2010s	107714560

From the above query , I have inferred that the results of decades per number of user votes with the voting that are extracted using the query using aggregation function .



From the above column chart , I have analyzed a insights of decades over the number of user votes.

## **Conclusion**

In Conclusion, I would like to conclude that IMDB Movie Analysis or any such analysis is done not only by Movie makers before movie production, but it is also done by various investors, stake holders, theatre outlet owners.

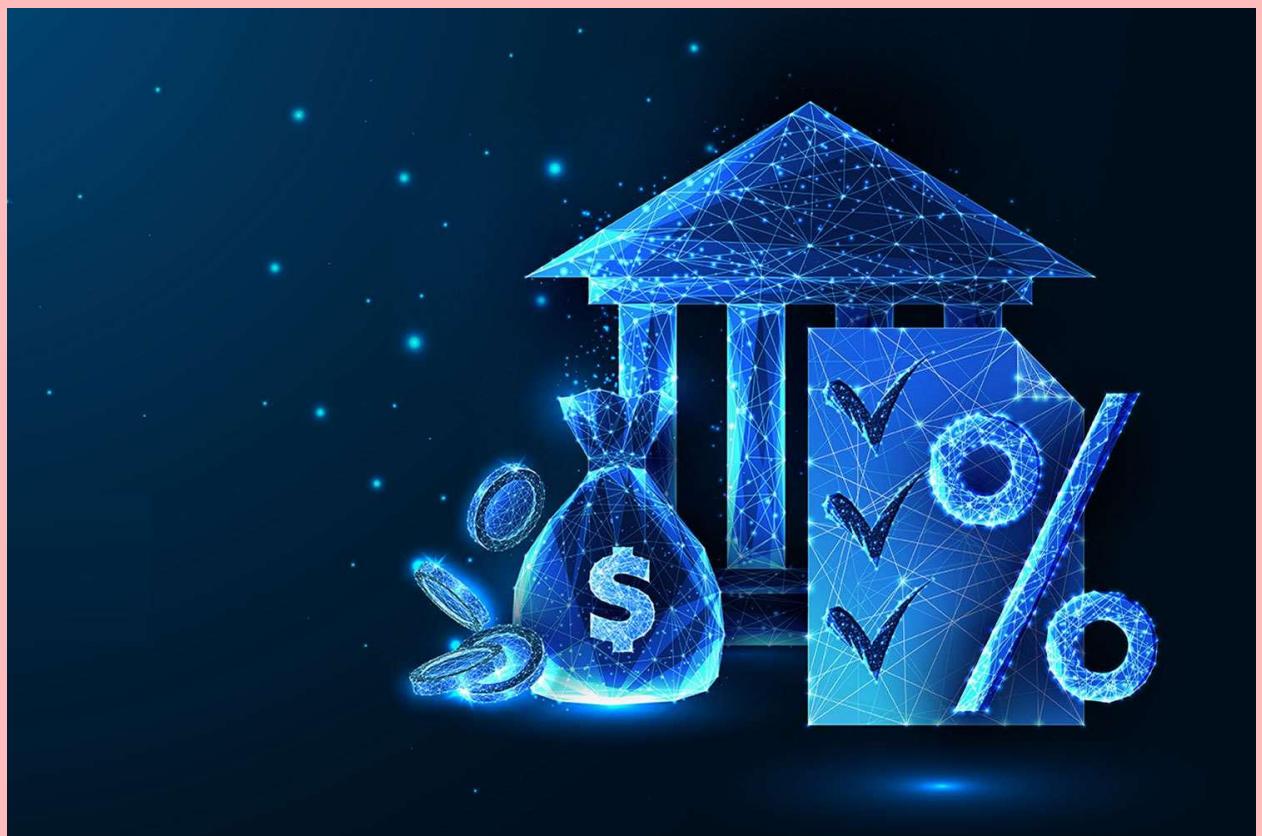
Normal people would not mind to do such analysis but such analysis plays an crucial part during the pre-production phase of the movies and also during the post-production phase.

Also, it is not necessary that the movie with the highest IMDB rating will have the highest profit. Profit is calculated truly on the basis on the number of tickets sold by theatres all over the world .

Most of the people are tired with their daily lives and they prefer movies with Comedy/ Drama genre or both, and they would not go for movies with Action/Horror genre So, directors and production team must keep in mind the above points and shall do the pre-production analysis before the commencement of filming.

# **Bank Loan Case Study**

## **Final Project-2**



## **DESCRIPTION**

Imagine you're a data analyst at a finance company that specializes in lending various types of loans to urban customers. Your company faces a challenge: some customers who don't have a sufficient credit history take advantage of this and default on their loans. Your task is to use Exploratory Data Analysis (EDA) to analyze patterns in the data and ensure that capable applicants are not rejected.

### **The Problem**

This case study aims to give you an idea of applying EDA in a real business scenario. In this case study, apart from applying the techniques that you have learnt in the EDA module, you will also develop a basic understanding of risk analytics in banking and financial services and understand how data is used to minimize the risk of losing money while lending to customers.

It aims to identify patterns which indicate if a client has difficulty paying their installments which may be used for taking actions such as denying the loan, reducing the amount of loan, lending (to risky applicants) at a higher interest rate, etc. This will ensure that the consumers capable of repaying the loan are not rejected.

Identification of such applicants using EDA is the aim of this case study. In other words, the company wants to understand the driving factors (or driver variables) behind loan default, i.e. the variables which are strong indicators of default. The company can utilize this knowledge for its portfolio and risk assessment.

To develop your understanding of the domain, you are advised to independently research a little about risk analytics – understanding the types of variables and their significance should be enough)

## Design

- First of all , in this project they have provided with a 2 datasets which are **current application** and **previous application**. So I have started cleaning the datasets and grouped as one, removed blanks, unwanted and unrelated columns and so on.
- I thought that in the given datasets unneeded columns are present there so I thought to remove the useless

columns which are not needed for risk analysis. I analyzed in given reference that helped me to analyze this risk assessment analysis.

- Next, I am supposed to find the outliers then removed the outliers, then I count the blanks of making it percentage how much the blanks are present so the percentage helped me to remove the unwanted columns and blanks by using excel **count functions**.
- By finding and cleaning all these data and moving on to tasks , my excel file named as **Final data** helped me to calculate the tasks and visually given a various insights like bar chart, box plot, column chart by make using excel pivot tables.

## Findings 1

A	B	C	D	E	F	G	H	I	J	K	L
6 of null values	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%
Count of rows	49999	49999	49999	49999	49999	49999	49999	49999	49999	49999	49998
SK_ID_CURR	TARGET	NAME_CONTRACT_TYPE	CODE_GENDER	FLAG_OWN_CAR	FLAG_OWN_REALTY	CNT_CHILDREN	AMT_INCOME_TOTAL	AMT_CREDIT	AMT_ANNUITY	AMT_GOOD	
100002	1	Cash loans	M	N	Y	0	202500	406597.5	24700.5		
100003	0	Cash loans	F	N	N	0	270000	1293502.5	35698.5		
100004	0	Revolving loans	M	Y	Y	0	67500	135000	6750		
100006	0	Cash loans	F	N	Y	0	135000	312682.5	29686.5		
100007	0	Cash loans	M	N	Y	0	121500	513000	21865.5		
100008	0	Cash loans	M	N	Y	0	99000	490495.5	27517.5		
100009	0	Cash loans	F	Y	Y	1	171000	1560726	41301		
100010	0	Cash loans	M	Y	Y	0	360000	1530000	42075		
100011	0	Cash loans	F	N	Y	0	112500	1019610	33826.5		
100012	0	Revolving loans	M	N	Y	0	135000	405000	20250		
100014	0	Cash loans	F	N	Y	1	112500	652500	21177		
100015	0	Cash loans	F	N	Y	0	38419.155	148365	10678.5		
100016	0	Cash loans	F	N	Y	0	67500	80865	5881.5		
100017	0	Cash loans	M	Y	N	1	225000	918468	28966.5		
100018	0	Cash loans	F	N	Y	0	189000	773680.5	32778		
100019	0	Cash loans	M	Y	Y	0	157500	299772	20160		
100020	0	Cash loans	M	N	N	0	108000	509602.5	26149.5		
100021	0	Revolving loans	F	N	Y	1	81000	270000	13500		
100022	0	Revolving loans	F	N	Y	0	112500	157500	7875		
100023	0	Cash loans	F	N	Y	1	90000	544491	17563.5		
100024	0	Revolving loans	M	Y	Y	0	135000	427500	21375		

A	B	C	D	E	F	G	H	I	J	K	L
DEX	SK_ID_CURR	TARGET	NAME_CONTRACT_TYPE	CODE_GENDER	FLAG_OWN_CAR	FLAG_OWN_REALTY	CNT_CHILDREN	AMT_INCOME_TOTAL	INCOME BIN	AMT	
1	100002	1	Cash loans	M	N	Y	0	202500	200K-225K		
2	100003	0	Cash loans	F	N	N	0	270000	250K-275K		
3	100004	0	Revolving loans	M	Y	Y	0	67500	50K-75K		
4	100006	0	Cash loans	F	N	Y	0	135000	125K-150K		
5	100007	0	Cash loans	M	N	Y	0	121500	100K-125K		
6	100008	0	Cash loans	M	N	Y	0	99000	75K-100K		
7	100009	0	Cash loans	F	Y	Y	1	171000	150K-175K		
8	100010	0	Cash loans	M	Y	Y	0	360000	350K-375K		
9	100011	0	Cash loans	F	N	Y	0	112500	100K-125K		
10	100012	0	Revolving loans	M	N	Y	0	135000	125K-150K		
11	100014	0	Cash loans	F	N	Y	1	112500	100K-125K		
12	100015	0	Cash loans	F	N	Y	0	38419.155	25K-50K		
13	100016	0	Cash loans	F	N	Y	0	67500	50K-75K		
14	100017	0	Cash loans	M	Y	N	1	225000	200K-225K		
15	100018	0	Cash loans	F	N	Y	0	189000	175K-200K		
16	100019	0	Cash loans	M	Y	Y	0	157500	150K-175K		
17	100020	0	Cash loans	M	N	N	0	108000	100K-125K		
18	100021	0	Revolving loans	F	N	Y	1	81000	75K-100K		
19	100022	0	Revolving loans	F	N	Y	0	112500	100K-125K		
20	100023	0	Cash loans	F	N	Y	1	90000	75K-100K		
21	100024	0	Revolving loans	M	Y	Y	0	135000	125K-150K		
22	100025	0	Cash loans	F	Y	Y	1	202500	200K-225K		
23	100026	0	Cash loans	F	N	N	1	450000	425K-450K		

From the above picture, I have represented the EDA analysis by percentage values and final dataset which are cleaned are pasted above as picture. Also then I have included a excel file as **EDA analysis** and **Final data**.

- SK\_ID\_PREV
- SK\_ID\_CURR
- NAME\_CONTRACT\_TYPE
- AMT\_ANNUITY

- AMT\_APPLICATION
- AMT\_CREDIT
- AMT\_GOODS\_PRICE
- WEEKDAY\_APPR\_PROCESS\_START
- HOUR\_APPR\_PROCESS\_START
- NAME\_CONTRACT\_STATUS
- DAYS\_DECISION
- NAME\_PAYMENT\_TYPE
- CODE\_REJECT\_REASON
- NAME\_CLIENT\_TYPE
- NAME\_GOODS\_CATEGORY
- CNT\_PAYMENT
- PRODUCT\_COMBINATION

A	B	C	D	E	F	G	H	I	J
% of null values	0%	0%	0%	21%	0%	0%	50%	21%	0%
Count of Rows	49999	49999	49999	39407	49999	49999	24801	39255	49999
SK_ID_PREV	SK_ID_CURR	NAME_CONTRACT_TYPE	AMT_ANNUITY	AMT_APPLICATION	AMT_CREDIT	AMT_DOWN_PAYMENT	AMT_GOODS_PRICE	WEEKDAY_APPR_PROCESS_START	H
2030495	271877	Consumer loans	1730.43	17145	17145	0	17145	SATURDAY	
2802425	108129	Cash loans	25188.615	607500	679671		607500	THURSDAY	
2523466	122040	Cash loans	15060.735	112500	136444.5		112500	TUESDAY	
2819243	176158	Cash loans	47041.335	450000	470790		450000	MONDAY	
1784265	202054	Cash loans	31924.395	337500	404055		337500	THURSDAY	
1383531	199383	Cash loans	23703.93	315000	340573.5		315000	SATURDAY	
2315218	175704	Cash loans		0	0			TUESDAY	
1656711	296299	Cash loans		0	0			MONDAY	
2367563	342292	Cash loans		0	0			MONDAY	
2579447	334349	Cash loans		0	0			SATURDAY	
1715995	447712	Cash loans	11368.62	270000	335754		270000	FRIDAY	
2257824	161140	Cash loans	13832.775	211500	246397.5		211500	FRIDAY	
2330894	258628	Cash loans	12165.21	148500	174361.5		148500	TUESDAY	
1397919	321676	Consumer loans	7654.86	53779.5	57564	0	53779.5	SUNDAY	
2273188	270658	Consumer loans	9644.22	26550	27252	0	26550	SATURDAY	
1232483	151612	Consumer loans	21307.455	126490.5	119853	12649.5	126490.5	TUESDAY	
2163253	154602	Consumer loans	4187.34	26955	27297	1350	26955	SATURDAY	
1285768	142748	Revolving loans	9000	180000	180000		180000	FRIDAY	
2393109	396305	Cash loans	10181.7	180000	180000		180000	THURSDAY	
1173070	199178	Cash loans	4666.5	45000	49455		45000	SATURDAY	
1506815	166490	Cash loans	25454.025	450000	491580		450000	MONDAY	

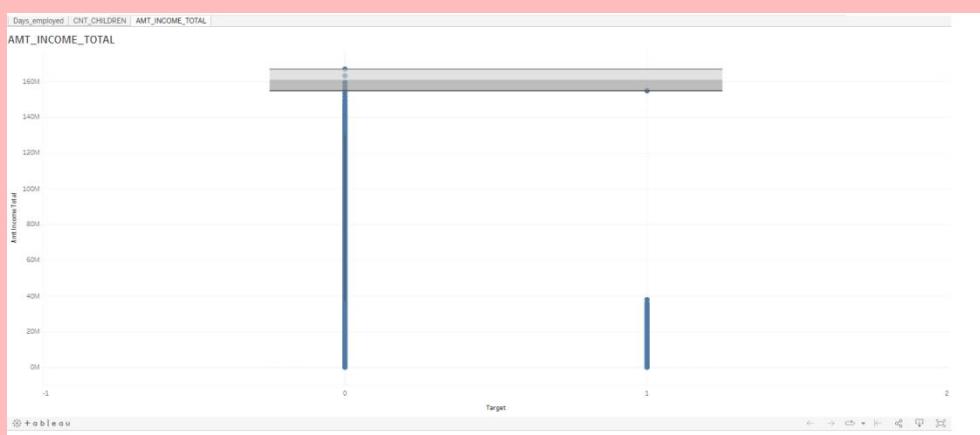
K_ID_PREV	SK_ID_CURR	NAME_CONTRACT_TYPE	AMT_ANNUITY	AMT_APPLICATION	AMT_CREDIT	AMT_GOODS_PRICE	WEEKDAY_APPR_PROCESS_START	HOUR_AP	NAME_CONTRACT_STATUS
2030495	271877	Consumer loans	1730.43	17145	17145	17145	SATURDAY		15 Approved
2802425	108129	Cash loans	25188.615	607500	679671	607500	THURSDAY		11 Approved
2523466	122040	Cash loans	15060.735	112500	136444.5	112500	TUESDAY		11 Approved
2819243	176158	Cash loans	47041.335	450000	470790	450000	MONDAY		7 Approved
1784265	202054	Cash loans	31924.395	337500	404055	337500	THURSDAY		9 Refused
1383531	199383	Cash loans	23703.93	315000	340573.5	315000	SATURDAY		8 Approved
1715995	447712	Cash loans	11368.62	270000	335754	270000	FRIDAY		7 Approved
2257824	161140	Cash loans	13832.775	211500	246397.5	211500	FRIDAY		10 Approved
2330894	258628	Cash loans	12165.21	148500	174361.5	148500	TUESDAY		15 Approved
1397919	321676	Consumer loans	7654.86	53779.5	57564	53779.5	SUNDAY		15 Approved
2273188	270658	Consumer loans	9644.22	26550	27252	26550	SATURDAY		10 Approved
1232483	151612	Consumer loans	21307.455	126490.5	119853	126490.5	TUESDAY		7 Approved
2163253	154602	Consumer loans	4187.34	26955	27297	26955	SATURDAY		12 Approved
1285768	142748	Revolving loans	9000	180000	180000	180000	FRIDAY		13 Approved
2393109	396305	Cash loans	10181.7	180000	180000	180000	THURSDAY		14 Approved
1173070	199178	Cash loans	4666.5	45000	49455	45000	SATURDAY		16 Refused
1506815	166490	Cash loans	25454.025	450000	491580	450000	MONDAY		6 Refused
1182516	267782	Cash loans	20361.6	405000	451777.5	405000	SATURDAY		4 Approved
1172937	302212	Cash loans	39475.305	1129500	1277104.5	1129500	THURSDAY		5 Refused
1543131	275707	Cash loans	22619.52	229500	241920	229500	THURSDAY		8 Approved
2536650	338725	Cash loans	16708.32	369000	369000	369000	WEDNESDAY		13 Approved
1676258	433469	Cash loans	22242.825	247500	268083	247500	THURSDAY		14 Approved
2075578	418383	Consumer loans	7656.705	74610	65610	74610	MONDAY		14 Approved

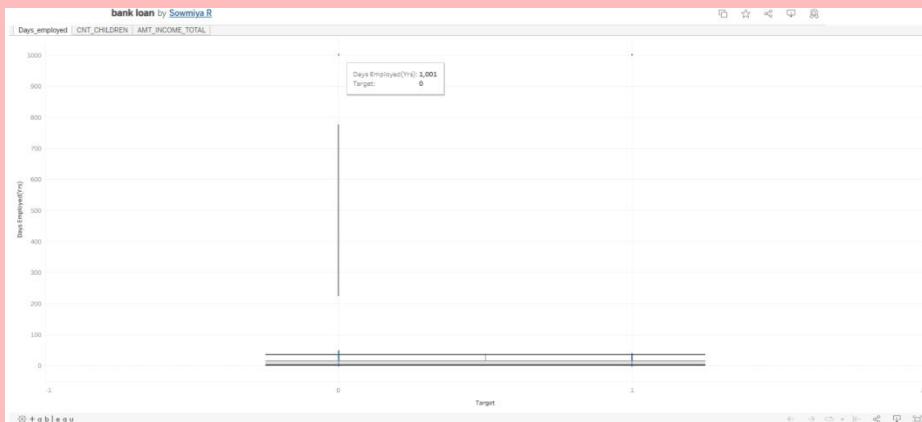
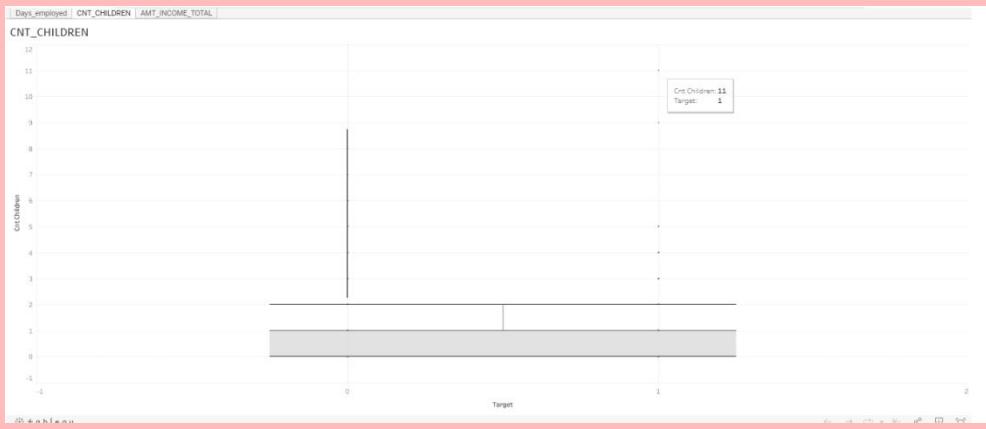
From the above 2 pictures , I had shown that counting the number of null rows with the percentage of every other columns and resulted to be 39256 entries and 17 header rows.

## Findings 2

Through then **AMT\_INCOME\_TOTAL**,

**CNT\_CHILDREN**, **DAYS\_EMPLOYED(yrs)** has some number of outliers.





I have used Tableau for finding out the outliers box plot which are visually represented in the above 3 images .

<b>Quartile 1</b>	87750
<b>Quartile 3</b>	248273
<b>Inter Quartile Range</b>	160523
<b>Upper Limit</b>	489057
<b>Lower Limit</b>	328534

This is the findings of outliers quartile ranges for amt\_income\_total which are shown above.



From the above dashboard is then outliers of amt\_credit, amt\_goods\_price, amt\_application, cnt\_payment, sk\_id\_curr, amt\_annuity and inferred that number of outliers with the relation between the previous application.

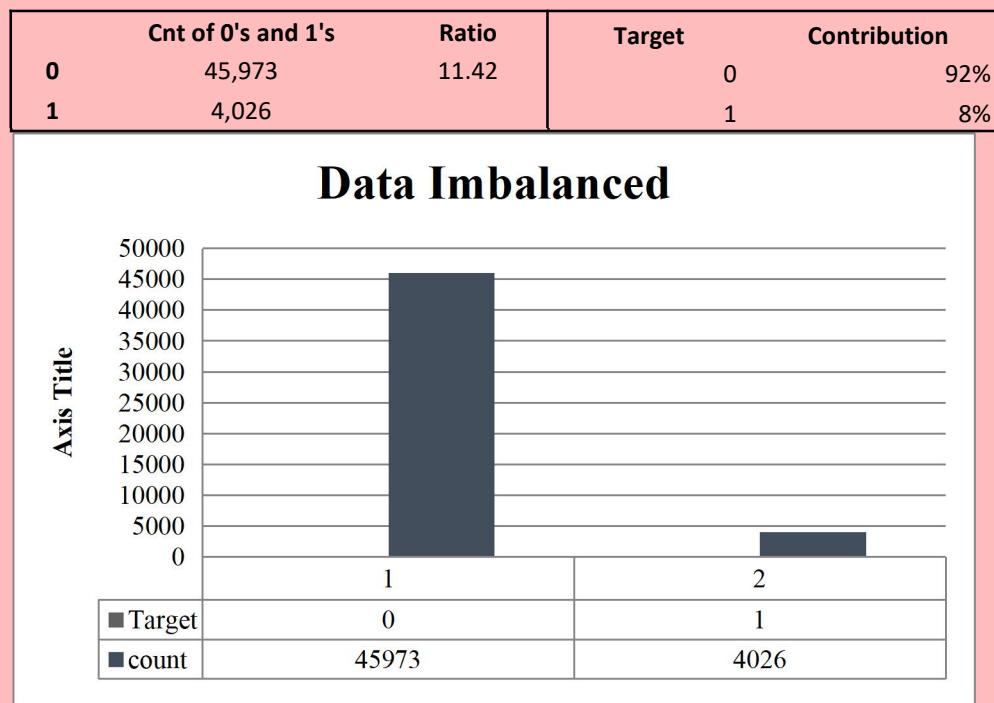
### Findings 3

In this data imbalance tasks, I have inferred that accuracy of data in the representation of visual charts as counting off the applicants which are grouped as **0** and **1** .

Target	count
0	45973
1	4026

<b>Grand Total</b>	49999
--------------------	-------

The above table says that the total number of applicants as Target (0 and 1) in the count which is total of 49999.



From the above chart represents the data imbalance of target as 0 and 1 (applicants) and the count of those referred to the repayment of loan status.

## Findings 4

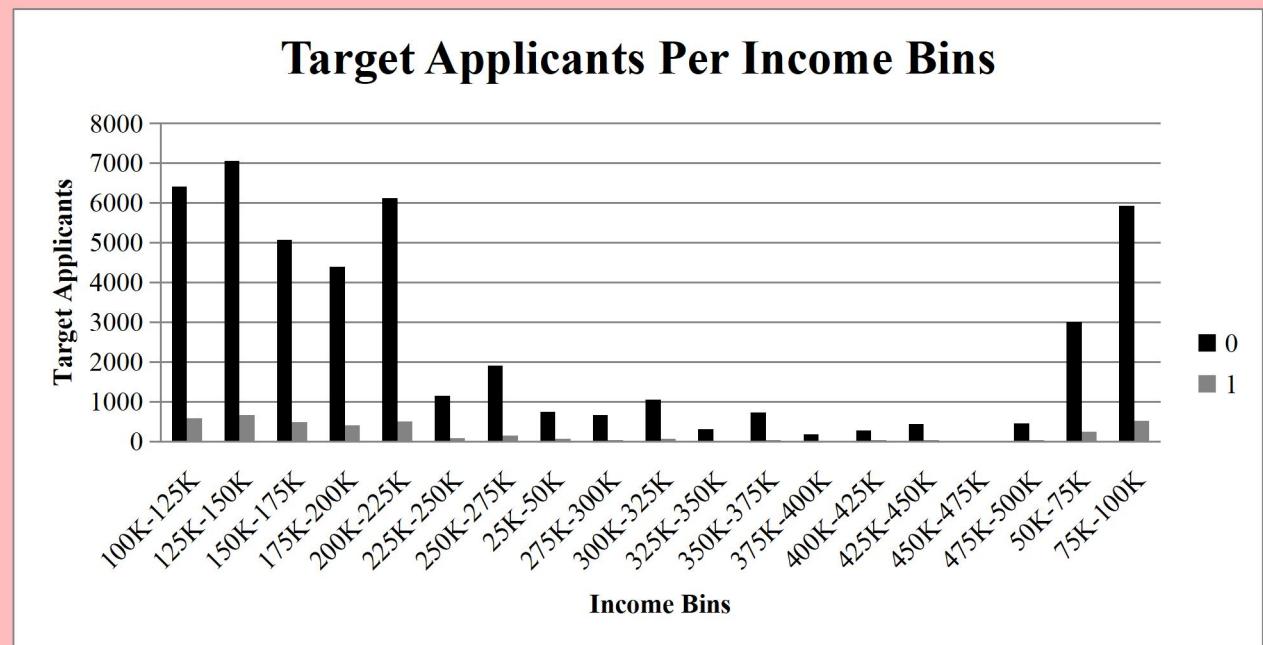
In this univariate and bivariate analysis , I had found out the class intervals of income bins and credit bins .

<b>Income</b>	<b>Income Bins</b>
0	0-25K
25,001	25K-50K
50,001	50K-75K
75,001	75K-100K
1,00,001	100K-125K
1,25,001	125K-150K
1,50,001	150K-175K
1,75,001	175K-200K
2,00,001	200K-225K
2,25,001	225K-250K
2,50,001	250K-275K
2,75,001	275K-300K
3,00,001	300K-325K
3,25,001	325K-350K
3,50,001	350K-375K
3,75,001	375K-400K
4,00,001	400K-425K
4,25,001	425K-450K
4,50,001	450K-475K
4,75,001	475K-500K
5,00,001	5 Lacs and above

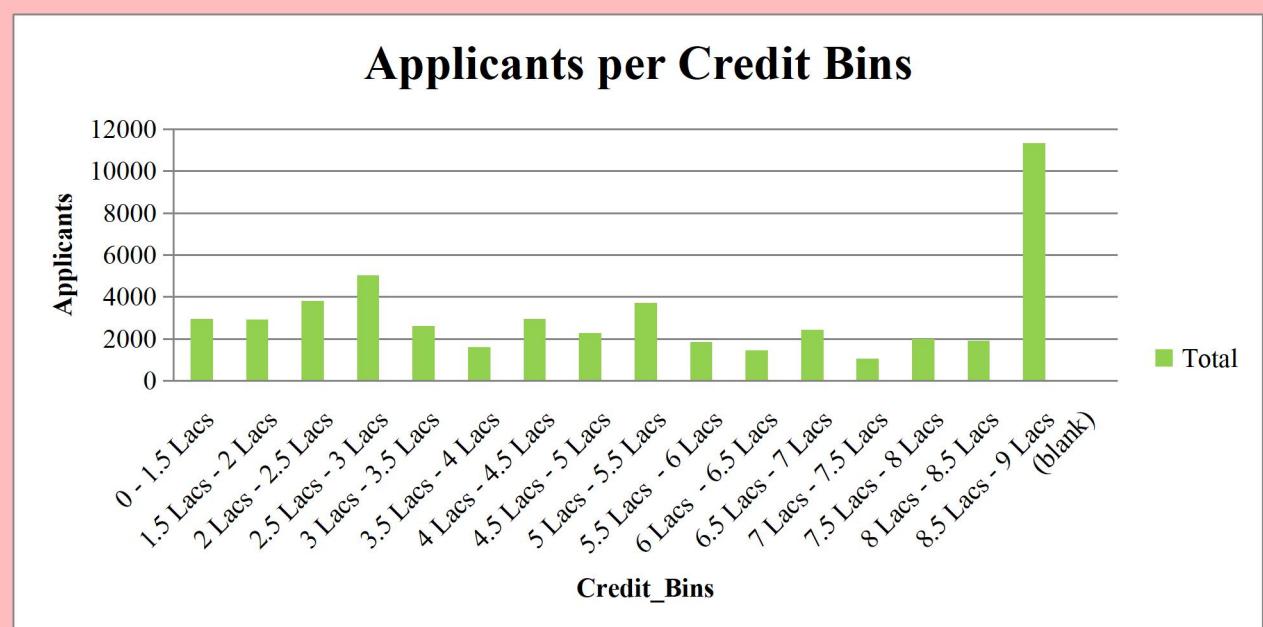
<b>Credit</b>	<b>Credit Bins</b>
0	0 - 1.5 Lacs
1,50,001	1.5 Lacs - 2 Lacs
2,00,001	2 Lacs - 2.5 Lacs
2,50,001	2.5 Lacs - 3 Lacs
3,00,001	3 Lacs - 3.5 Lacs
3,50,001	3.5 Lacs - 4 Lacs
4,00,001	4 Lacs - 4.5 Lacs
4,50,001	4.5 Lacs - 5 Lacs
5,00,001	5 Lacs - 5.5 Lacs
5,50,001	5.5 Lacs - 6 Lacs
6,00,001	6 Lacs - 6.5 Lacs
6,50,001	6.5 Lacs - 7 Lacs
7,00,001	7 Lacs - 7.5 Lacs
7,50,001	7.5 Lacs - 8 Lacs
8,00,001	8 Lacs - 8.5 Lacs
8,50,001	8.5 Lacs - 9 Lacs
9,00,001	9 Lacs and above

These are the class intervals which are to found in the univariate , segmented univariate, bivariate analysis.

## SEGMENTED UNIVARIATE

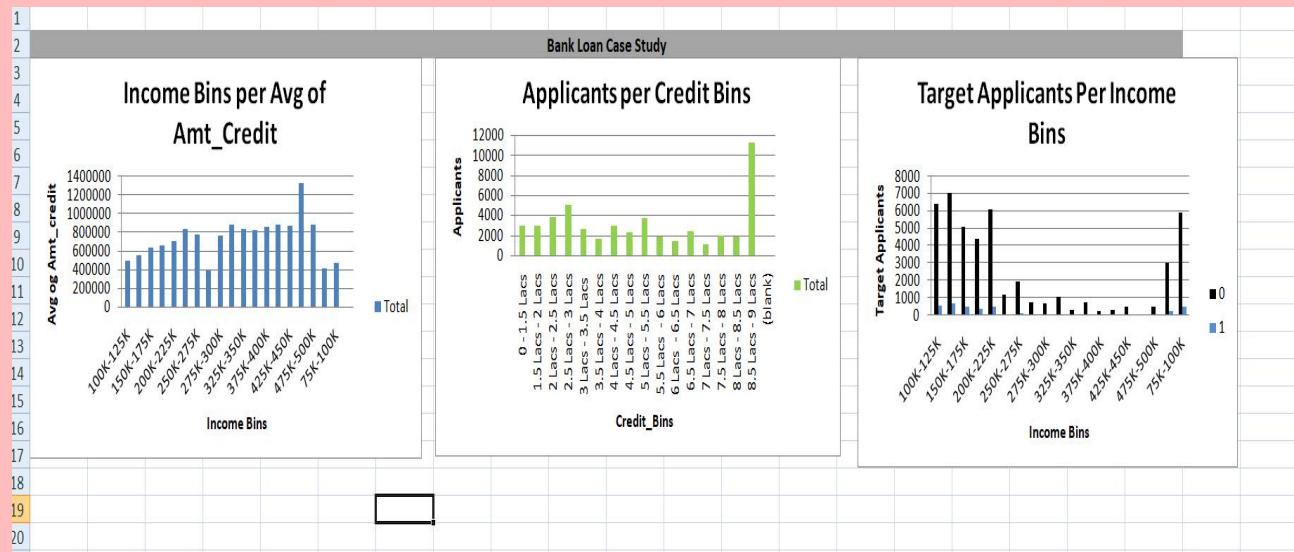
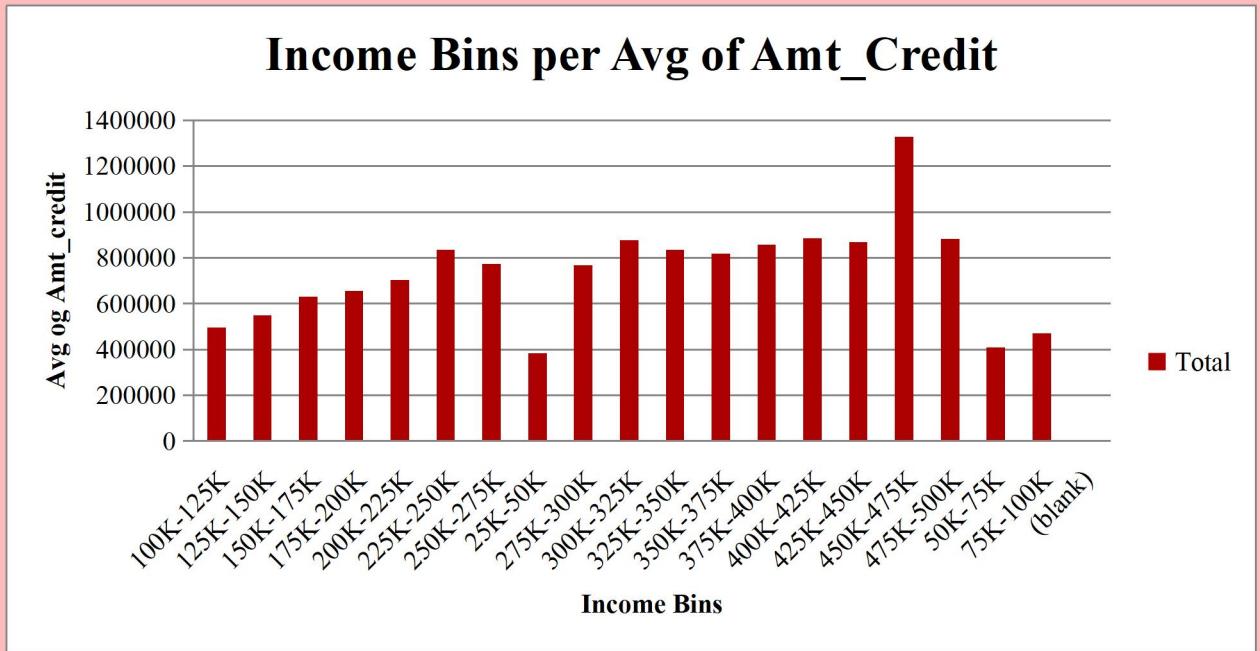


## UNIVARIATE ANALYSIS



From the above chart , I analyzed that the credit bins is the score of applicants pay the loan or not.

## BIVARIATE ANALYSIS



## Findings 5

- i) So I have taken up of **target 0** as one correlation. Then I included these columns alone from the final dataset and performed correlation between them.

**TARGET, CNT\_CHILDREN, AMT\_INCOME\_TOTAL, AMT\_CREDIT,**

**REGION\_POPULATION\_RELATIVE, DAYS\_BIRTH(yrs),**

**DAYS\_EMPLOYED(YRS), DAYS\_ID\_PUBLISH(YRS),**

**REGION\_RATING\_CLIENT**

		CORRELATION FOR APPLICANTS WITH PAYMENT MADE ON TIME							
		1	0.047	0.011	-0.026	-0.322	-0.250	0.044	0.011
CNT_CHILDREN	1	0.047	0.011	-0.026	-0.322	-0.250	0.044	0.011	
AMT_INCOME_TOTAL	0.047	1	0.406	0.175	-0.073	-0.183	-0.033	-0.231	
AMT_CREDIT	0.011	0.406	1	0.070	0.052	-0.083	0.019	-0.103	
REGION_POPULATION_RELATIVE	-0.026	0.175	-0.026	1	0.033	-7E-03	-0.001	-0.534	
DAYS_BIRTH(YEARS)	-0.322	-0.073	0.052	0.033	1	0.633	0.262	0.003	
DAYS_EMPLOYED(YRS)	-0.250	-0.183	-0.083	7E-05	0.633	1	0.259	0.043	
DAYS_ID_PUBLISH(YRS)	0.044	-0.033	0.019	-0.001	0.262	0.259	1	0.015	
REGION_RATING_CLIENT	0.011	-0.231	-0.103	-0.534	0.003	0.043	0.015	1	
		CNT_CHILDREN	AMT_INCOME_TOTAL	AMT_CREDIT	REGION_POPULATION_RELATIVE	DAYS_BIRTH(YEARS)	DAYS_EMPLOYED(YRS)	DAYS_ID_PUBLISH(YRS)	REGION_RATING_CLIENT

So from the above selected columns I have correlated and inferred to create a heat map for this **Target 0**. In the above heat map itself says the answer that is the correlation of each and every different scenarios.

- ii) For the **Target 1** correlation, the same as above the calculations as including the columns of particular needed headed rows and performing calculation of correlation of different scenarios.

TARGET, CNT\_CHILDREN, AMT\_INCOME\_TOTAL, AMT\_CREDIT,  
 REGION\_POPULATION\_RELATIVE, DAYS\_BIRTH(yrs),  
 DAYS\_EMPLOYED(YRS), DAYS\_ID\_PUBLISH(YRS),  
 REGION\_RATING\_CLIENT

With this columns I have inferred a **correlation of target 1** in the kind of heat maps.

CORRELATION FOR APPLICANTS WITH PAYMENT DIFFICULTIES									
CNT_CHILDREN	1	-0.068	0.053	-0.009	-0.235	-0.162	0.100	-0.025	
AMT_INCOME_TOTAL	-0.068	1	0.379	0.144	0.039	-0.108	-0.019	-0.144	
AMT_CREDIT	0.053	0.379	1	0.061	0.165	-0.043	0.095	-0.015	
REGION_POPULATION_RELATIVE	-0.009	0.144	0.061	1	-0.052	-0.114	0.022	-0.498	
DAYS_BIRTH(Years)	-0.235	0.039	0.165	-0.052	1	0.545	0.288	0.100	
DAYS_EMPLOYED (Years)	-0.162	-0.108	-0.043	-0.114	0.545	1	0.224	0.090	
DAYS_ID_PUBLISH(Years)	0.100	-0.019	0.095	0.022	0.288	0.224	1	0.019	
REGION_RATING_CLIENT	-0.025	-0.144	-0.015	-0.498	0.100	0.090	0.019	1	
CNT_CHILDREN AMT_INCOME_TOTAL AMT_CREDIT REGION_POPULATION_RELATIVE DAYS_BIRTH(Years) DAYS_EMPLOYED (Years) DAYS_ID_PUBLISH(Years) REGION_RATING_CLIENT									

From the above heat map , I have analyzed that the correlation between the different scenarios of **correlation for target 1** and above map itself calculated for this tasks for correlation of each and every other loan scenarios.

## Conclusion

- The proportion/percentage of the defaulters(target = 1) is around 8% and that of non-defaulters(target = 0) is around 92% The Bank generally lends more loan to

Female clients as compared to Males clients as the count of Female clients in the defaulter's list is less than that of Males. Still Bank can look for more Male clients if their credit amount is satisfied.

- Also the clients who belong to Working class tend to pay their loans on time followed by the clients who fall under Commercial Associate
- Clients having Education status like Secondary/ Higher Secondary or more tend to pay loan on time so bank can prefer lending loans to clients having such Education Status Clients who fall in the Age Group 31-40 have the highest count for paying off their loans on time followed by the clients who fall in the Age Groups 41-60
- Clients having LOW credit amount range tend to pay off their loans on time than compared to HIGH and MEDIUM credit range

- Clients living with their Parents tend to pay off their loans quickly as compared to other housing type. So Bank can lend loan to clients having housing type → Living with Parents
- Clients taking loan for purchasing New Home i.e. clients taking Home Loans or purchasing New Car i.e. Car Loans and clients who have a income type as State Servant tend to pay their loans on time and hence Bank should prefer clients having such background
- The Bank should be more cautious when lending money to clients with Repairs purpose because they have high count of Defaulters along with High count of Defaulters

# Analyzing the Impact of Car Features on Price and Profitability

Final Project-3



## **DESCRIPTION**

The automotive industry has been rapidly evolving over the past few decades, with a growing focus on fuel efficiency, environmental sustainability, and technological innovation. With increasing competition among manufacturers and a changing consumer landscape, it has become more important than ever to understand the factors that drive consumer demand for cars.

In recent years, there has been a growing trend towards electric and hybrid vehicles and increased interest in alternative fuel sources such as hydrogen and natural gas. At the same time, traditional gasoline-powered cars remain dominant in the market, with varying fuel types and grades available to consumers.

For the given dataset, as a Data Analyst, the client has asked How can a car manufacturer optimize pricing and product development decisions to maximize profitability while meeting consumer demand?

This problem could be approached by analyzing the relationship between a car's features, market category, and pricing, and identifying which features and categories are most popular among consumers and most profitable for the manufacturer. By using data analysis techniques such as regression analysis and market segmentation, the manufacturer could develop a pricing strategy that balances

consumer demand with profitability, and identify which product features to focus on in future product development efforts. This could help the manufacturer improve its competitiveness in the market and increase its profitability over time.

## **The Problem**

### **TASKS: ANALYSIS**

1. How does the popularity of a car model vary across different market categories?
2. What is the relationship between a car's engine power and its price?
3. Which car features are most important in determining a car's price?
4. How does the average price of a car vary across different manufacturers?
5. What is the relationship between fuel efficiency and the number of cylinders in a car's engine?

## **BUILDING THE DASHBOARD:**

1. How does the distribution of car prices vary by brand and body style?
2. Which car brands have the highest and lowest average MSRPs, and how does this vary by body style?
3. How do the different features such as transmission type affect the MSRP, and how does this vary by body style?
4. How does the fuel efficiency of cars vary across different body styles and model years?
5. How does the car's horsepower, MPG, and price vary across different Brands?

## **Design**

Perform statistical analysis to understand the distribution, variability, and basic characteristics of the collected data. Visualize data patterns, correlations, and anomalies using graphs and plots. Identify initial insights that can guide subsequent analysis. Evaluate model performance using metrics like accuracy, precision, recall, and F1-score. Showcase how predictive models can optimize maintenance schedules and reduce downtime.

Demonstrate how real-time alerts and warnings can be generated for drivers and relevant authorities. Evaluate the potential reduction in accidents and overall enhancement in road safety. Analyze driving patterns and vehicle performance data to uncover factors influencing fuel efficiency. Develop a model to predict fuel efficiency based on driving behavior, road conditions, and vehicle characteristics.

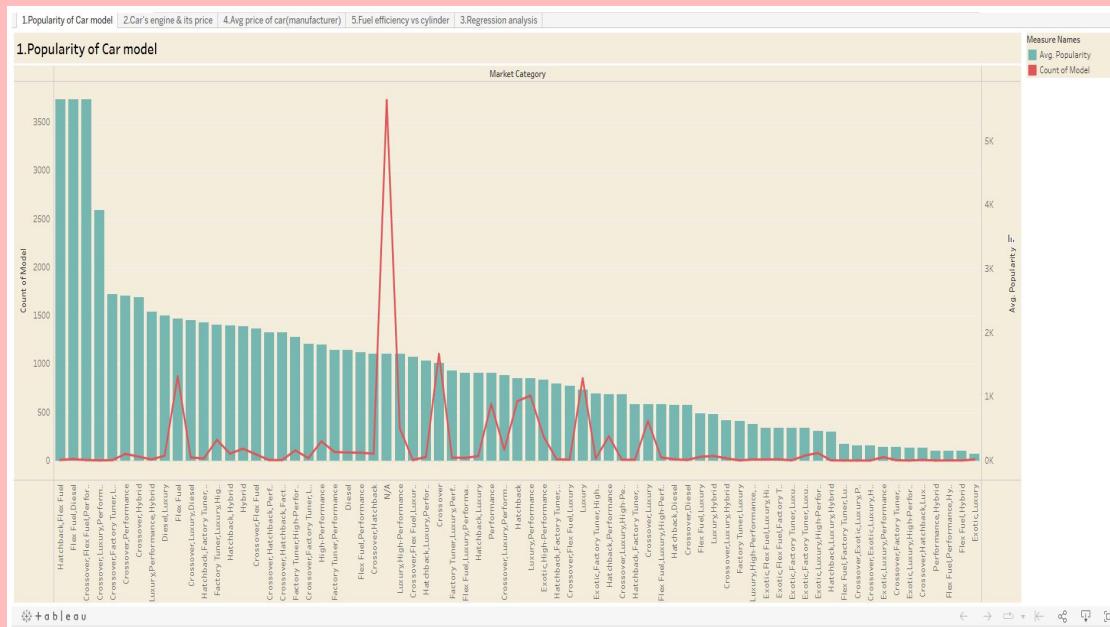
The project aims to provide a comprehensive understanding of the impact of car data analysis on the automotive industry, offering insights, solutions, and guidance for industry stakeholders to harness the power of data for positive transformations.

## Findings 1

The solution for the task 1 of A is that of created a pivot table with the help of showing the number of car models in each market category with respect to the popularity scores which is of creating a make or manufacturers with the popularity scores of maintaining the market category of impact of car trends.

Popularity of a car model vary across different market categories		
Market Category	Values	
	Average of Popularity	Count of Model
Crossover	1545	1110
Crossover,Diesel	873	7
Crossover,Exotic,Luxury,High-Performance	238	1
Crossover,Exotic,Luxury,Performance	238	1
Crossover,Factory Tuner,Luxury,High-Performance	1823	26
Crossover,Factory Tuner,Luxury,Performance	2607	5
Crossover,Factory Tuner,Performance	210	4
Crossover,Flex Fuel	2074	64
Crossover,Flex Fuel,Luxury	1173	10
Crossover,Flex Fuel,Luxury,Performance	1624	6
Crossover,Flex Fuel,Performance	5657	6
Crossover,Hatchback	1676	72
Crossover,Hatchback,Factory Tuner,Performance	2009	6
Crossover,Hatchback,Luxury	204	7
Crossover,Hatchback,Performance	2009	6
Crossover,Hybrid	2563	42
Crossover,Luxury	885	410
Crossover,Luxury,Diesel	2149	34
Crossover,Luxury,High-Performance	1037	9
Crossover,Luxury,Hybrid	631	24
Crossover,Luxury,Performance	1345	113
Crossover,Luxury,Performance,Hybrid	3916	2
Crossover,Performance	2586	69
Diesel	1731	84
Diesel,Luxury	2275	51
Exotic,Factory Tuner,High-Performance	1046	21
Exotic,Factory Tuner,Luxury,High-Performance	518	52

Exotic,Factory Tuner,Luxury,Performance	520	3
Exotic,Flex Fuel,Factory Tuner,Luxury,High-Performance	520	13
Exotic,Flex Fuel,Luxury,High-Performance	520	11
Exotic,High-Performance	1271	261
Exotic,Luxury	113	12
Exotic,Luxury,High-Performance	467	79
Exotic,Luxury,High-Performance,Hybrid	204	1
Exotic,Luxury,Performance	217	36
Exotic,Performance	1391	10
Factory Tuner,High-Performance	1941	106
Factory Tuner,Luxury	617	2
Factory Tuner,Luxury,High-Performance	2133	215
Factory Tuner,Luxury,Performance	1413	31
Factory Tuner,Performance	1696	92
Flex Fuel	2217	872
Flex Fuel,Diesel	5657	16
Flex Fuel,Factory Tuner,Luxury,High-Performance	258	1
Flex Fuel,Hybrid	155	2
Flex Fuel,Luxury	747	39
Flex Fuel,Luxury,High-Performance	879	33
Flex Fuel,Luxury,Performance	1380	28
Flex Fuel,Performance	1680	87
Flex Fuel,Performance,Hybrid	155	2
Hatchback	1319	641
Hatchback,Diesel	873	14
Hatchback,Factory Tuner,High-Performance	1205	13
Hatchback,Factory Tuner,Luxury,Performance	887	9
Hatchback,Factory Tuner,Performance	2159	22
Hatchback,Flex Fuel	5657	7
Hatchback,Hybrid	2121	72
Hatchback,Luxury	1380	46
Hatchback,Luxury,Hybrid	454	3
Hatchback,Luxury,Performance	1566	38
Hatchback,Performance	1040	252
High-Performance	1821	199
Hybrid	2106	123
Luxury	1103	855
Luxury,High-Performance	1668	334
Luxury,High-Performance,Hybrid	569	12
Luxury,Hybrid	674	52
Luxury,Performance	1293	673
Luxury,Performance,Hybrid	2333	11
Performance	1349	601
Performance,Hybrid	155	1
<b>Grand Total</b>	<b>1499</b>	<b>8172</b>



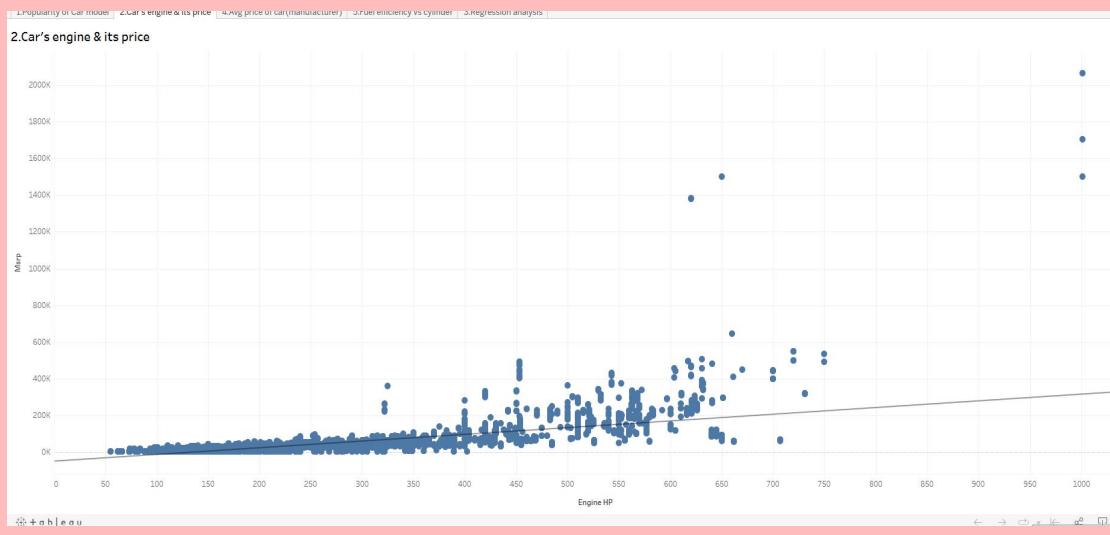
From the above combo chart tasks, I have created with a bar chart and combination of line chart, it varies as colors for both line and bar chart. Here I have insisted that the popularity of car is based on the average of popularity and count of each model.

## Findings 2

the relationship between car's engine and its price as MSRP , found out by separating from the final dataset of cleaned data which help me out to categorize the particular tasks and have been visualized with the trend line chart.

	A	B	C	D
1	Engine HP	MSRP		
2	335	\$ 46,135		
3	300	\$ 40,650		
4	300	\$ 36,350		
5	230	\$ 29,450		
6	230	\$ 34,500		
7	230	\$ 31,200		
8	300	\$ 44,100		
9	300	\$ 39,300		
10	230	\$ 36,900		
11	230	\$ 37,200		
12	300	\$ 39,600		
13	230	\$ 31,500		
14	300	\$ 44,400		
15	230	\$ 37,200		
16	230	\$ 31,500		
17	320	\$ 48,250		
18	320	\$ 43,550		
19	172	\$ 2,000		
20	172	\$ 2,000		
21	172	\$ 2,000		
22	172	\$ 2,000		
23	172	\$ 2,000		
24	172	\$ 2,000		

The above picture depicts that the separated data of car's **Engine HP** and **MSRP**, to find the relationship between the car's engine power and its price.



From the above picture depicts that the scatter plot and combination of shown trend lines for finding the relationship between Engine HP and MSRP with respect to it.

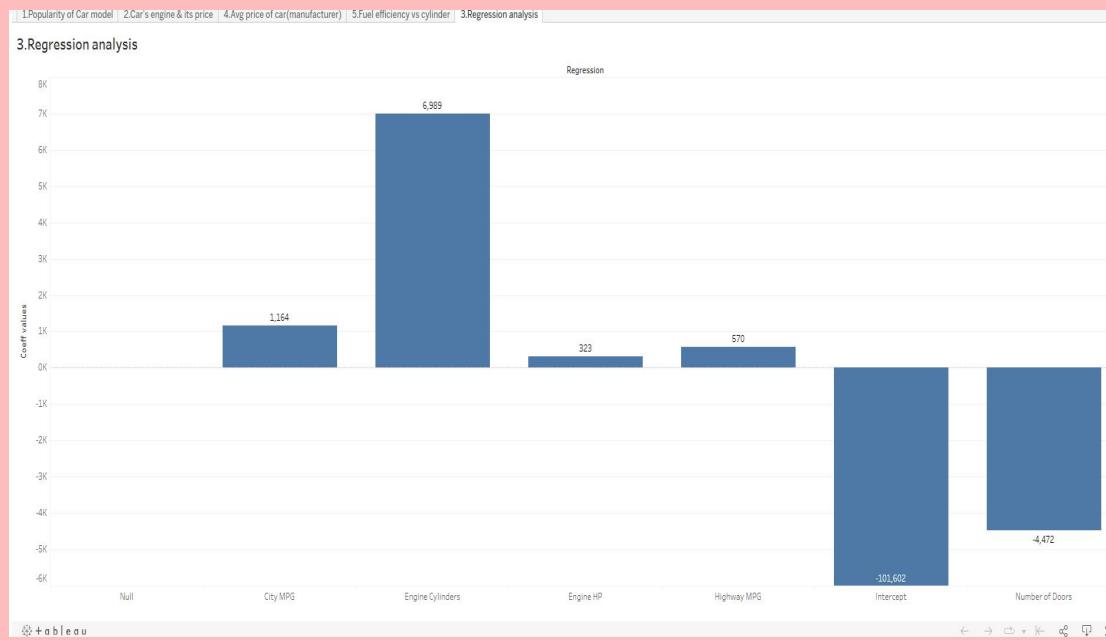
Then I have taken for finding the visualization of scatter plot and with a trend lines is then taken Engine HP and MSRP to solve the relationship between them.

And also Engine HP will be taken as x-axis and MSRP taken as y-axis and plotting it accordingly and gone through the relation between , in the analysis it will be shown as trend lines - show trend lines and then edit the axis according to it.

## Findings 3

Here I have learned how to solve the regression analysis using the excel data analysis functions, it gives the anova table , summary output with the residual and regression which inferred the regression data of car features on determining a car's price.

SUMMARY OUTPUT								
Regression Statistics								
Multiple R	0.680708139							
R Square	0.46336357							
Adjusted R Square	0.463136297							
Standard Error	44170.77827							
Observations	11812							
ANOVA								
	df	SS	MS	F	Significance F			
Regression	5	1.98891E+13	3.97782E+12	2038.799457	0			
Residual	11806	2.30342E+13	1951057653					
Total	11811	4.29233E+13						
	Coefficients	Standard Error	t Stat	P-value	Lower 95%	Upper 95%	Lower 95.0%	Upper 95.0%
Intercept	-101601.736	3684.351697	-27.57655738	2.765E-162	-108823.673	-94379.79896	-108823.673	-94379.79896
Engine HP	322.7465574	6.01767382	53.63310924	0	310.9509241	334.5421906	310.9509241	334.5421906
Engine Cylinders	6989.177662	439.6449924	15.89732121	2.53591E-56	6127.400961	7850.954363	6127.400961	7850.954363
Number of Doors	-4472.158125	465.7180593	-9.602715711	9.35015E-22	-5385.042338	-3559.273912	-5385.042338	-3559.273912
Highway MPG	570.1808088	105.7839778	5.390048859	7.17937E-08	362.826764	777.5348535	362.826764	777.5348535
City MPG	1163.755457	121.9978136	9.539150113	1.72109E-21	924.61962	1402.891294	924.61962	1402.891294

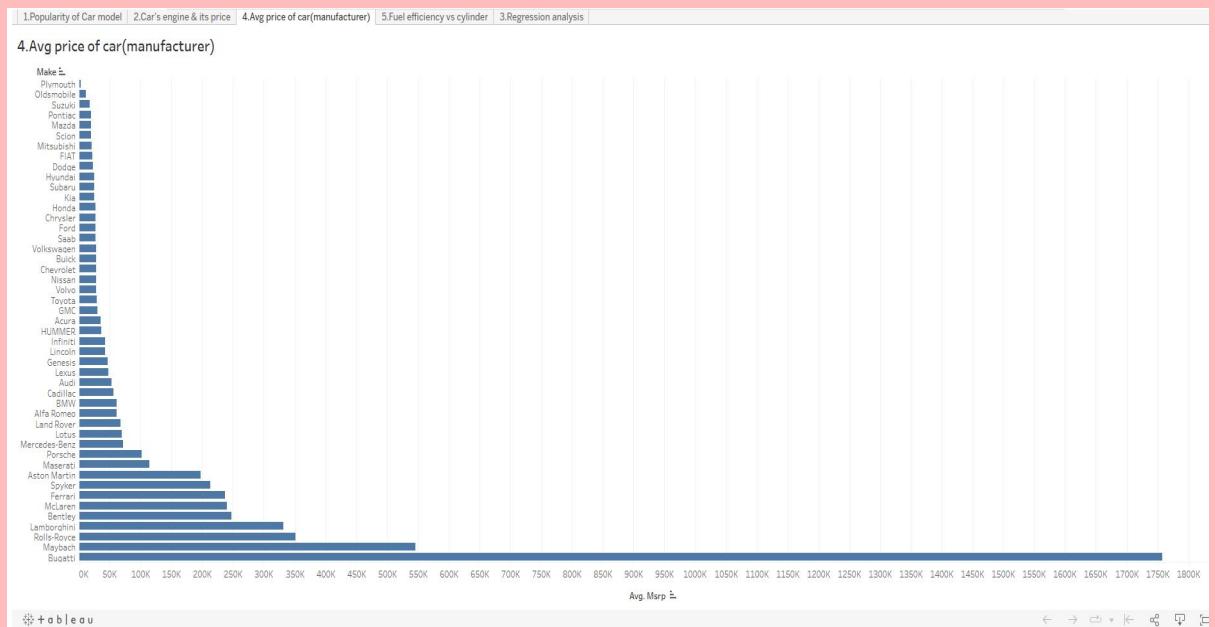


## **Findings 4**

So I have insisted on creating a pivot table with the help of manufacturers and MSRP for price of each make. And then created the pivot table with the row labels of make that is manufacturer and MSRP as Average of car price, that are included as pivot table and then insisted as the changing of count to average of car price with the help of value field settings in the pivot table.

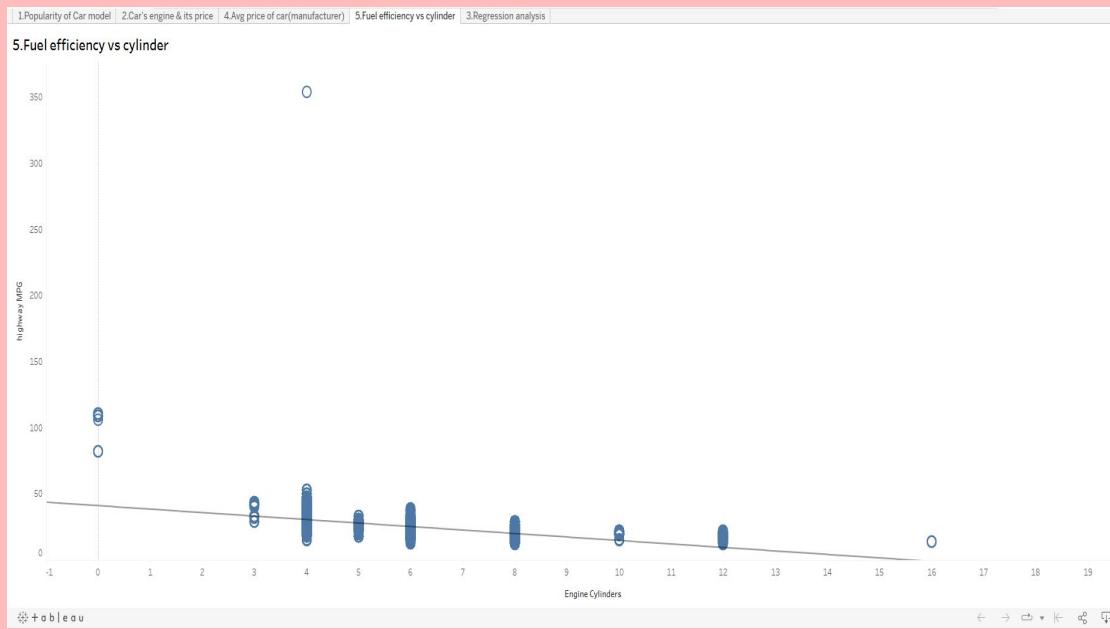
<b>Manufacturer</b>	<b>Average of MSRP</b>
Acura	₹ 34,887.59
Alfa Romeo	₹ 61,600.00
Aston Martin	₹ 1,97,910.38
Audi	₹ 53,452.11
Bentley	₹ 2,47,169.32
BMW	₹ 61,546.76
Bugatti	₹ 17,57,223.67
Buick	₹ 28,206.61
Cadillac	₹ 56,231.32
Chevrolet	₹ 28,273.36
Chrysler	₹ 26,722.96
Dodge	₹ 22,390.06
Ferrari	₹ 2,37,383.82
FIAT	₹ 22,206.02
Ford	₹ 27,393.42
Genesis	₹ 46,616.67
GMC	₹ 30,493.30
Honda	₹ 26,629.82
HUMMER	₹ 36,464.41
Hyundai	₹ 24,597.04
Infiniti	₹ 42,394.21
Kia	₹ 25,112.39
Lamborghini	₹ 3,31,567.31
Land Rover	₹ 67,823.22
Lexus	₹ 47,549.07
Lincoln	₹ 42,494.37
Lotus	₹ 69,188.28
Maserati	₹ 1,14,207.71
Maybach	₹ 5,46,221.88
Mazda	₹ 19,719.06
McLaren	₹ 2,39,805.00

Mercedes-Benz	₹	71,537.81
Mitsubishi	₹	21,215.47
Nissan	₹	28,513.37
Oldsmobile	₹	11,542.54
Plymouth	₹	3,122.90
Pontiac	₹	19,321.55
Porsche	₹	1,01,622.40
Rolls-Royce	₹	3,51,130.65
Saab	₹	27,413.50
Scion	₹	19,932.50
Spyker	₹	2,13,323.33
Subaru	₹	24,827.50
Suzuki	₹	17,900.96
Toyota	₹	28,946.15
Volkswagen	₹	28,076.20
Volvo	₹	28,541.16
<b>Grand Total</b>	₹	<b>40,559.94</b>



From the above pasted picture represents that the horizontal a stacked bar chart for insights of average price of car with respect to the manufacturer. Taking Average of price that is in the header row of MSRP is been taken as x-axis and Manufacturer will be taken as y-axis .

## Findings 5

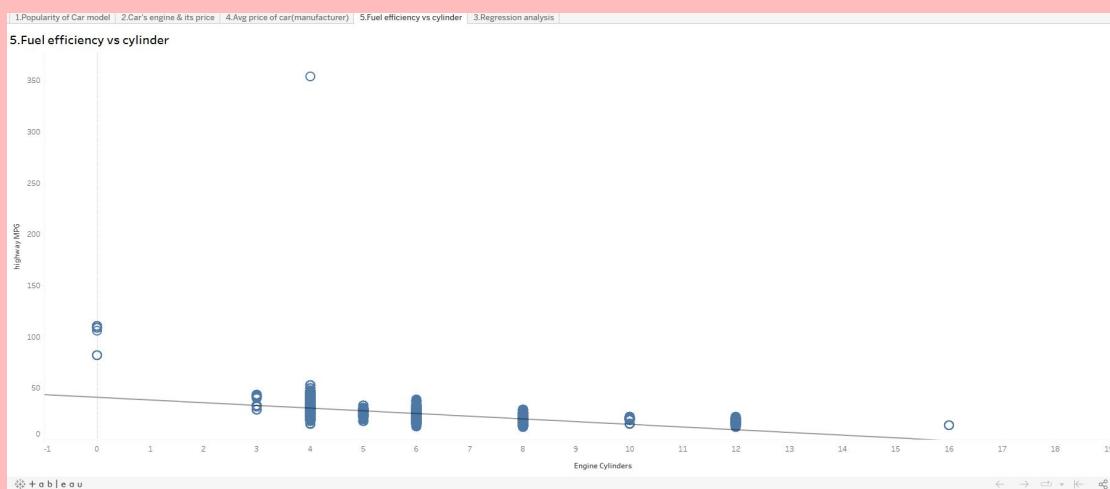


The above picture represents the scatter plot of finding the fuel efficiency versus number of cylinder and with a trend lines for the reference of insisting the correlation of fuel efficiency with each of the cylinders.

I found the solution that the correlation of coefficient between number of cylinders and the highway MPG played a role in this tasks. It quantifies the strength of each of the cylinders numbered and a direction of relationship between them.

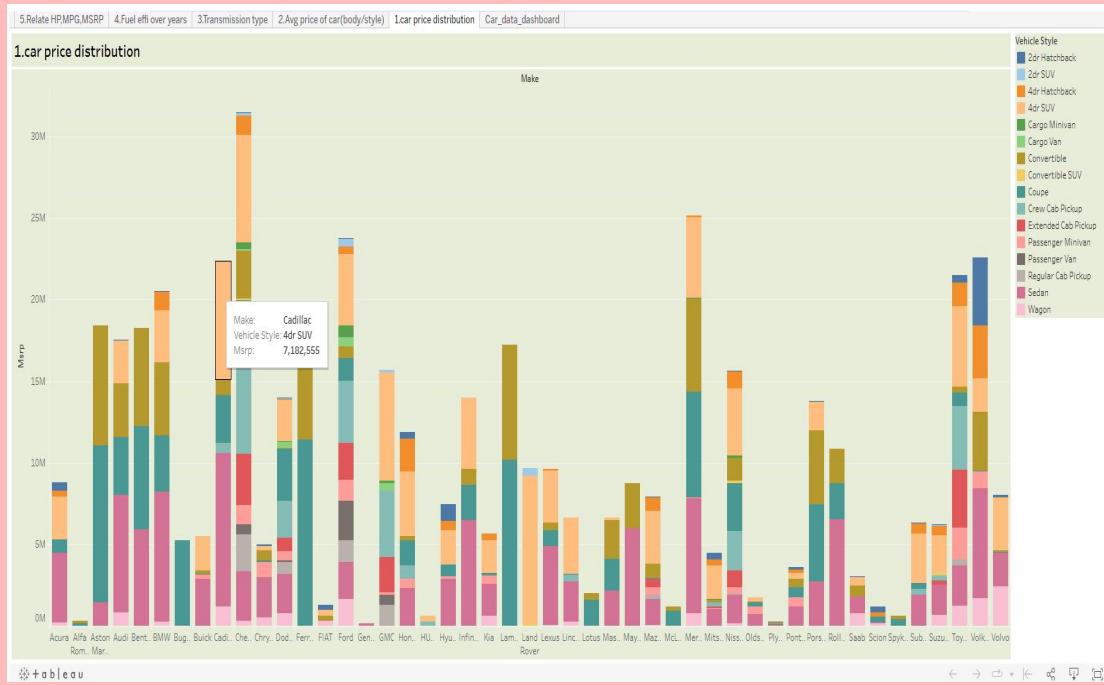
	A	B	C	D	E	F
1	Engine Cylinders	highway MPG				
2	6	26			Correlation, r=	-0.62031
3	6	28				
4	6	28				
5	6	28				
6	6	28				
7	6	28				
8	6	26				
9	6	28				
10	6	28				
11	6	27				
12	6	28				
13	6	28				
14	6	28				
15	6	28				
16	6	28				
17	6	25				
18	6	28				
19	6	24				
20	6	24				
21	6	20				
22	6	24				
23	6	21				
24	6	24				

From the above picture is my analysis of finding the fuel efficiency with the number of cylinders and found out that coefficient values as **r=-0.620312551.**



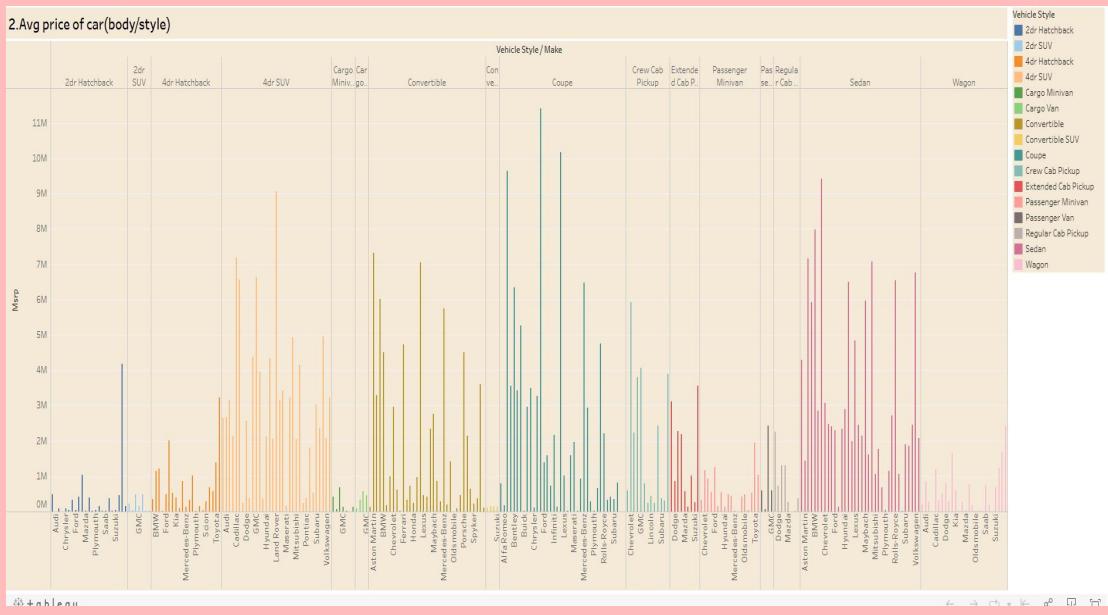
# Building dashboards

## Findings 1



In this dashboard analysis, I have inferred that the distribution of car price with respect to the vehicle style and manufacturer for each of the car. Now here what is the tasks is that to make a stacked column chart for making the visualizations of finding the distribution of car price based on branch and vehicle style.

## Findings 2



For this dashboard analysis , I have insisted on finding the average price car based on both body and style position. So I have created a column chart to visualize the values here I have done both the body and style together with each of the brand with its price that is MSRP.

## Findings 3



For this dashboard analysis, I have inferred that the transmission type of vehicle styles as a feature that affects the MSRP of car to shown below the pivot tables. Here I have taken vehicle style as x-axis and msrp as y-axis to get the designed scatter plot to visualize the relationship between the msrp and vehicle style, and also legends shown as for the reference of each of the color given to vehicle styles.

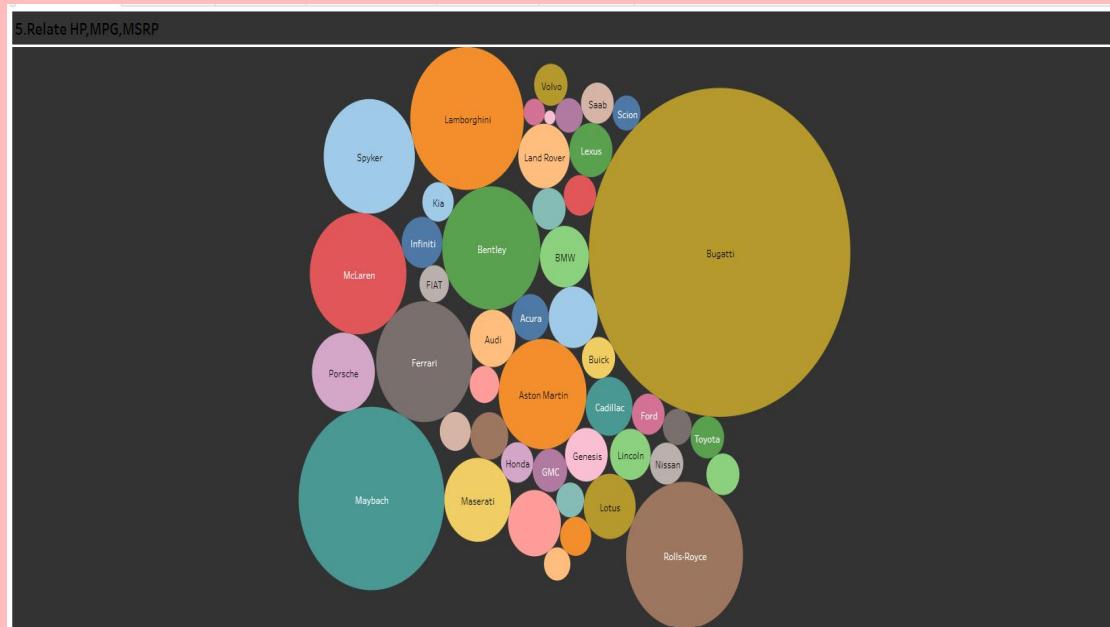
## **Findings 4**

For this dashboard analysis, it is make a Line chart to show the trend of fuel efficiency (MPG) over time for each body style. It represents that the fuel efficiency of finding the average number of highway MPG with respect to the year.

And calculated the average MPG for each combination of body style and model year with the visualization of line chart with the trend lines.



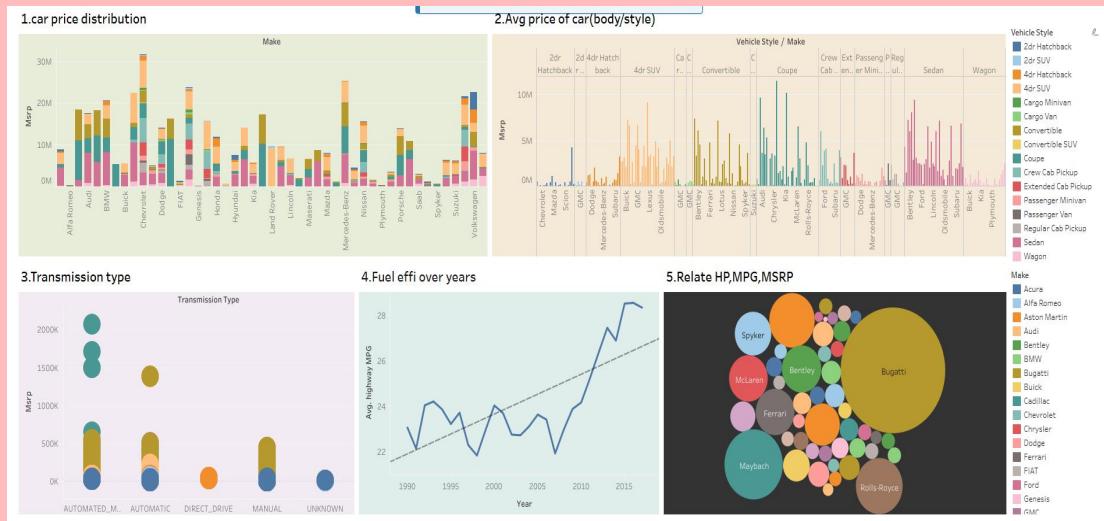
## **Findings 5**



For this dashboard analysis, I have insisted on making a bubble chart with the help of relationship between horsepower, MPG, and price across different car brands. Assign different colors to each brand and label the bubbles

with the car model name. And have calculated the average horsepower, MPG, and MSRP for each car brand.

## DASHBOARD(overall)



## Conclusion

In this project , I have learnt how to make a analysis of each of the dashboards, I was very thrilled to making a beautiful insights as a data analyst. I have used tableau public and excel to find the tasks and dashboard analysis.

Future directions continue to be improvements in software and algorithms. Increasing the level of parallelism is desired for dealing with large scale computational and memory requirements.

# ABC CALL VOLUME TREND ANALYSIS

Final Project4



## **DESCRIPTION**

In this project, you'll be diving into the world of 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 such as the agent's name and ID, the queue time (how long a customer had to wait before connecting with an agent), the time of the call, the duration of the call, and the call status (whether it was abandoned, answered, or transferred).

A Customer Experience (CX) team plays a crucial role in a company. They analyze customer feedback and data, derive insights from it, and share these insights with the rest of the organization. This team is responsible for a wide range of tasks, including managing customer experience programs, handling internal communications, mapping customer journeys, and managing customer data, among others.

In the current era, several AI-powered tools are being used to enhance customer experience. These include Interactive Voice Response (IVR), Robotic Process Automation (RPA), Predictive Analytics, and Intelligent Routing.

One of the key roles in a CX team is that of the customer service representative, also known as a call center agent. These agents handle various types of support, including email, inbound, outbound, and social media support.

## **The Problem**

Calculate the average call time duration for all incoming calls received by agents (in each Time\_Bucket).

Show the total volume/ number of calls coming in via charts/ graphs [Number of calls v/s Time]. You can select time in a bucket form (i.e. 1-2, 2-3, .....) )

As you can see current abandon rate is approximately 30%. Propose a manpower plan required during each time bucket [between 9am to 9pm] to reduce the abandon rate to 10%. (i.e. You have to calculate minimum number of agents required in each time bucket so that at least 90 calls should be answered out of 100.)

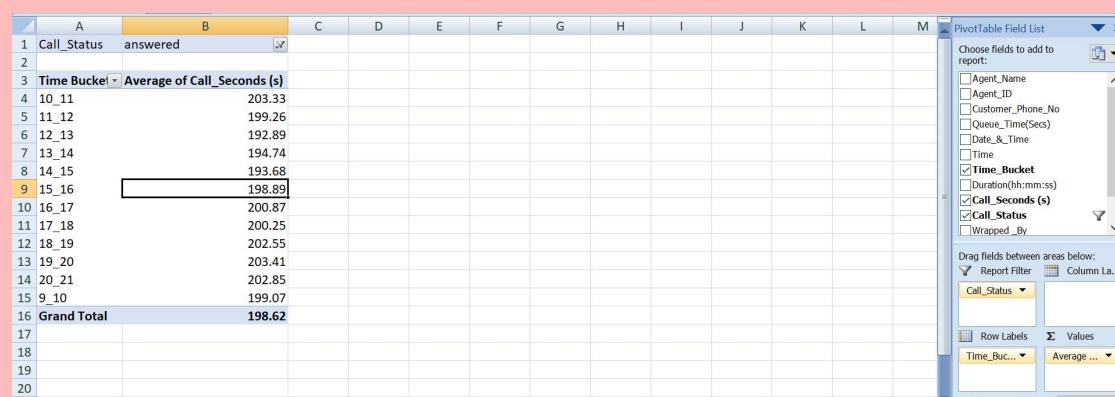
Suppose every 100 calls that customer made during 9 Am to 9 Pm, customer also made 30 calls in night between interval [9 Pm to 9 Am] and distribution of those 30 calls are as follows: Now propose a manpower plan required during each time bucket in a day. Maximum Abandon rate assumption would be same 10%.

Assumption: An agent work for 6 days a week; On an average total unplanned leaves per agent is 4 days a month; An agent total working hrs is 9 Hrs out of which 1.5 Hrs goes into lunch and snacks in the office. On average an agent occupied for 60% of his total actual working Hrs (i.e 60% of 7.5 Hrs) on call with customers/ users. Total days in a month is 30 days.

## Findings 1

For this tasks , that is, average duration of calls for each time bucket have found out to be average of call duration for each of the time bucket. So I have been analyzed in the cleaned data and created a pivot table based on the understanding of the question.

Call_Status	answered
Time Bucket	Average of Call_Seconds (s)
10_11	203.33
11_12	199.26
12_13	192.89
13_14	194.74
14_15	193.68
15_16	198.89
16_17	200.87
17_18	200.25
18_19	202.55
19_20	203.41
20_21	202.85
9_10	199.07
<b>Grand Total</b>	<b>198.62</b>



From the above excel picture represents that the pivot table calculation for finding the average of duration of calls with

each of the time bucket and analyzed with pivot table calculation of finding the average of call duration.

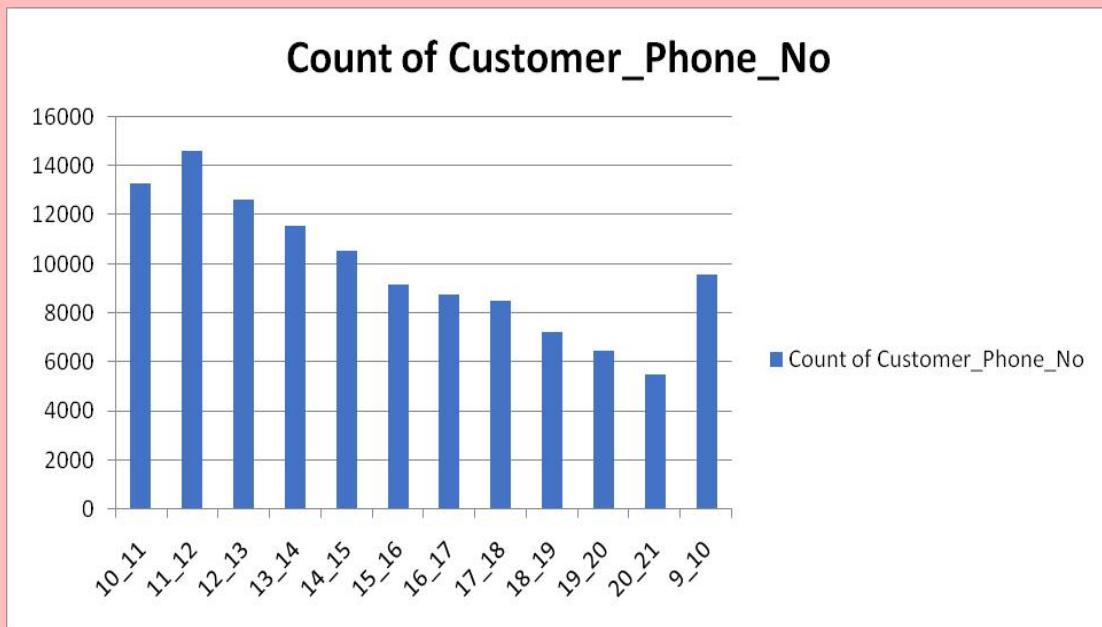
So with that , the average of call duration for each time bucket is **198.62**

## Findings 2

For this tasks, I created the pivot table about the number of calls received with each of the time bucket, so taken with a time bucket, number of time and customer phone number.

Values		
Time Bucket	Count of Customer_Phone_No	Count of Time
10_11	13313	11%
11_12	14626	12%
12_13	12652	11%
13_14	11561	10%
14_15	10561	9%
15_16	9159	8%
16_17	8788	7%
17_18	8534	7%
18_19	7238	6%
19_20	6463	5%
20_21	5505	5%
9_10	9588	8%
<b>Grand Total</b>	<b>117988</b>	<b>100%</b>

From the above pivot table chart represents that count of time and with each of the customer phone number with the time bucket of 10\_11 to 9\_10.



The above bar chart depicts that the time bucket with each of the customer phone number also that grand total of **117988** with count of time would be 100%.

### Findings 3

Distribution of 30 calls coming in night for every 100 calls coming in between 9am - 9pm (i.e. 12 hrs slot)											
9pm- 10pm	10pm - 11pm	11pm- 12am	12am- 1am	1am - 2am	2am - 3am	3am - 4am	4am - 5am	5am - 6am	6am - 7am	7am - 8am	8am - 9am
3	3	2	2	1	1	1	1	1	3	4	4

From this assumptions, I have been analyzed through the simple mathematical operations in the microsoft excel, then delivered a distribution of time buckets of 30 calls coming in night for every 100 calls.

107018	23-01-2022 20:58		1	1		
107019	23-01-2022 20:58		1	1		
107020	23-01-2022 20:59		1	1		
107021	<b>Grand Total</b>	<b>34403</b>	<b>82452</b>	<b>1133</b>	<b>117988</b>	
107022		1496	3585	49	5130	
107023		29%	70%	1%	100%	
107024						
107025	Working Hrs / agent	4.5				
107026	Avg call duration	198.62				
107027						
107028						
107029	For 90% hrs needed	254.73015				
107030	No of agents need is	57				
107031			9am-9pm		57	90%
107032			9pm-9am		17	90%
107033	Avg no of calls in night is	1539		overall agents needed is	74	
107034	To increase call rate to 90% ni8 is	76				
107035	Num of agents needed in night	17				

So above picture is the calculation for finding the minimum number of agents required in each time bucket . Then I have calculated with the pivot table calculations and picked a **abandon (total value/no of values)**

**Answered(total value/no of values)**

**Transfer(total value/no of values)**

With the count of duration(hh:mm:ss).

Working Hrs / agent	4.5
Avg call duration	198.62

This is calculated for finding the working hours per each of the agent is 4.5 and Average of call duration is 198.62(previous calculation). This is been calculated as ...

$$\text{Working hours per agent} = (60/100) * 7.5 \\ = 4.5$$

Average of call duration=198.62

Time Bucket	Average of Call_Seconds (s)
-------------	--------------------------------

10_11	203.33
11_12	199.26
12_13	192.89
13_14	194.74
14_15	193.68
15_16	198.89
16_17	200.87
17_18	200.25
18_19	202.55
19_20	203.41
20_21	202.85
9_10	199.07
<b>Grand Total</b>	<b>198.62</b>

For 90% hrs needed	254.73015
No of agents need is	57

From this we can calculate the number of agents needed for 90Hrs is 57. This can be calculated as .....

$$\text{For 90hrs needed} = 5130 * 198.62 * 0.9 / 3600 \\ = 254.73015$$

$$\text{No of agents needed} = 254.73015 / 4.5 \\ = 57$$

## Findings 4

For this tasks , I have used previously calculated pivot table and used with respect to the analysis of tasks, This creates a poor customer experience. Assume that for every 100 calls that customers make between 9 am and 9 pm, they also make 30 calls at night between 9 pm and 9 am. So I have done calculations based on the night shift between 9pm-9am.

Avg no of calls in night is	1539
To increase call rate to 90% ni8 is	76
Num of agents needed in night	17

In this excel sheet calculation, I have analyzed with a average number of calls in night is 1539 and to increase the call rate to 90% in night is 76 and number of agents have calculated for the night shift.

The calculations be....

$$\text{Average number of calls in night} = 0.3 * 5130$$

$$= 1539$$

$$\begin{aligned} \text{To increase the call rate to 90\% in} \\ \text{night} &= 198.62 * 1539 * 0.9 / 3600 \end{aligned}$$

$$= 76$$

$$\text{Number of agents needed in night} = 76 / 4.5$$

$$= 17$$

9am-9pm	57	90%
9pm-9am	17	90%
overall agents needed is	74	

Then the overall agents is been listed above and the calculations would be....

$$9am-9pm = 57 \text{ (90\%)}$$

$$9pm-9am = 17 \text{ (90\%)}$$

$$\text{And the overall agents needed} = 57 + 17$$

$$= 74$$

## **Conclusion**

In this project ABC Call Volume Trend Analysis we analysis customer experience. A customer experience team consists of professionals who analyse customer feedback and data, and share insights with the rest of the organization. Keep a good customer relationship that help business for future growth. The project done through using Microsoft excel. So it help me to increase my technical skills and knowledge in excel. In this analysis project I have analysis the customer relationship, to solve customer problem that keep good relationship between the business and customer. After doing this project this help me to improve my data analytical skills, visualization skills etc.

## **Appendix**

### **1. Data Analytics Process:-**

---> Link for the shared PDF on Google Drive:

[https://drive.google.com/file/d/1KMCuiG0fkwXOrVK9fDC1cKvg3u-zLhUJ/view?usp=drive\\_link](https://drive.google.com/file/d/1KMCuiG0fkwXOrVK9fDC1cKvg3u-zLhUJ/view?usp=drive_link)

### **2. Instagram User Analytics:-**

----> Link for the shared file on Google Drive :

[https://drive.google.com/file/d/10-vCx1VXrsek946TH1n2bWMbeyWHhGEV/view?usp=drive\\_link](https://drive.google.com/file/d/10-vCx1VXrsek946TH1n2bWMbeyWHhGEV/view?usp=drive_link)

### **3. Operation Analytics and Investigating Metric Spike Analysis:-**

-----> Link for the shared file on Google Drive:

[https://drive.google.com/file/d/1Bs5mmol7THWJmUfB-hHOZ5MKXZ-e9n0a/view?usp=drive\\_link](https://drive.google.com/file/d/1Bs5mmol7THWJmUfB-hHOZ5MKXZ-e9n0a/view?usp=drive_link)

### **4. Hiring Process Analytics:-**

----> Link for shared PDF on google drive:

[https://drive.google.com/file/d/13hIf2XhhMIQtwI70n5BKBuNt8afX5DaD/view?usp=drive\\_link](https://drive.google.com/file/d/13hIf2XhhMIQtwI70n5BKBuNt8afX5DaD/view?usp=drive_link)

## **5. IMDB Movie Analysis-**

**---> Link for the shared PDF on Google Drive:**

[https://drive.google.com/drive/folders/1L1Lmf1OXD2d6N2DZx08B4saOG2RrdUwi?usp=drive\\_link](https://drive.google.com/drive/folders/1L1Lmf1OXD2d6N2DZx08B4saOG2RrdUwi?usp=drive_link)

## **6. Bank Loan Case Study:-**

**----> Link for the shared file on Google Drive:**

[https://drive.google.com/drive/folders/1BPXUmps2cb8Q4bD1HmBHuwSYyl5?usp=drive\\_link](https://drive.google.com/drive/folders/1BPXUmps2cb8Q4bD1HmBHuwSYyl5?usp=drive_link)

## **7. Analyzing the Impact of Car Features on Price and Profitability:-**

**----> Link for the shared file on Google Drive:**

[https://drive.google.com/drive/folders/1teT0DcMzWfXiVx\\_33PFkRoUtyh\\_v9gJK?usp=drive\\_link](https://drive.google.com/drive/folders/1teT0DcMzWfXiVx_33PFkRoUtyh_v9gJK?usp=drive_link)

## **8. ABC Call Volume Trend Analysis:-**

**-----> Link for the shared file on Google Drive:**

[https://drive.google.com/drive/folders/19utlO7Ypv0pzRLC-4Q3vXB-lw4zpOQU0?usp=drive\\_link](https://drive.google.com/drive/folders/19utlO7Ypv0pzRLC-4Q3vXB-lw4zpOQU0?usp=drive_link)

**Link to GitHub Portfolio:-**

<https://github.com/Sowmiyar1512>

**Link to LinkedIn Profile:-**

<https://www.linkedin.com/in/sowmiya-r-58767a215/>