

DATA ANALYSIS PORTFOLIO

*PREPARED BY
HUMAIRA KHAN*

PROFESSIONAL BACKGROUND

I am a third-year cybersecurity student pursuing bachelors as my degree from the University of Pune. I have secured 9.0 CGPA and am well-versed with Python, SQL, NoSQL, Statistics, Excel, Machine Learning and Data Analysis.

I have worked on several projects such as steganography using python, breast cancer prediction using SVM algorithm and movie recommendation system.

Coming to my data analysis projects, I have analyzed census data, analyzed NASA planetary budget exploration using SQL, uber trips analysis.

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 have theoretical knowledge. But I am waiting to use my theoretical knowledge in a practical way. And I believe by putting significant efforts, I will learn.



Table Of Contents

Professional Background -----	2
Table of Contents -----	3
Data Analytics Process -----	4
Instagram User Analytics-----	5-13
Operation Analytics and Investigating Metric Spike-----	14-25
Hiring Process Analytics -----	26-31
IMDB Movie Analysis-----	32-47
Bank Loan Case Study-----	48-63
XYZ Ads Airing Report Analysis-----	64-77
ABC Call Volume Trend Analysis-----	78-86
Conclusion-----	87

DATA ANALYTICS PROCESS ORDERING FOOD FROM ONLINE APPS

PLAN: SO, FIRST WE WILL DECIDE FROM WHAT APP WE WILL ORDER FOOD FROM, IS IT ZOMATO OR SWIGGY? WHAT WE FEEL LIKE EATING, INDIAN OR CHINESE? SWEET OR SPICY?

PREPARE: NEXT I NEED TO HOW MUCH AM I WILLING TO SPEND AND HOW TO GET THAT MONEY.

PROCESS: THEN I NEED TO CHECK HOW MUCH I WANT FROM THE DATA. IF I WANT TO EAT INDIAN FOOD, WHAT EXACTLY AND HOW MUCH AMOUNT I WANT. 2 CHAPATIS WITH SABZI, WILL IT BE ENOUGH?

ANALYZE: BUT WAIT. I AM DIET CONSCIOUS. WILL MY ORDER SYNCHRONIZE WITH MY CALORIES INTAKE? WILL IT PROVIDE ME WITH THE PERFECT AMOUNT OF NUTRIENTS?

SHARE: NOW THAT YOU HAVE ANALYSED AND DECIDED ON THE PERFECT ORDER, YOU WOULD WANT TO SHARE YOUR IDEA BY CLICKING ON WHAT YOU WANT TO BUY.

ACT: AND THEN YOU FINALLY ORDER THE FOOD AND LATER EAT IT.

INSTAGRAM USER ANALYTICS USING SQL

PROJECT DESCRIPTION:

THIS PROJECT IS ABOUT ANALYSING HOW USERS ENGAGE WITH A VERY POPULAR SOCIAL APP NAMED INSTAGRAM. WHICH USER HAS THE MOST NUMBER OF LIKES AND FOLLOWERS? WHICH USER IS BARELY ACTIVE? HOW MUCH TIME DOES AN AVERAGE USER SPENDS DAILY ON INSTAGRAM? WHICH DAY ARE PEOPLE MOST ACTIVE? TO ANSWER ALL THESE CURIOUS QUESTIONS, WE WILL BE USING SQL AND ANALYSE EVERYTHING ABOUT IT BY THE INSTAGRAM DATA WE ALREADY HAVE.

APPROACH:

AFTER INSTALLATION AND EXECUTING THE DATABASE NAMED IG_CLONE, I FIRST USED THE SHOW TABLE COMMANDS AND THEN THE SHOW COLUMNS COMMAND TO UNDERSTAND ABOUT ALL THE TABLE, COLUMNS AND THE DATA INSIDE IT. THEN READING THE SQL QUERY QUESTIONS AND UNDERSTANDING WHAT ALL DATA WE EXACTLY NEED, I STARTED QUERIES ABOUT THE SAME.

INSIGHTS:

LETS START WITH THE QUERIES IN ORDER AND WHAT ALL WE LEARNED BY IT.

MARKETING QUERIES

THE 5 OLDEST USERS.

The screenshot shows two instances of MySQL Workbench running side-by-side. Both instances have the same interface setup, including the Navigator pane on the left, the Query Editor pane in the center, and the SQLAdditions pane on the right.

Session 1 (Top):

- Query:** `SELECT * FROM users ORDER BY created_at LIMIT 5`
- Result Grid:**

ID	username	created_at
1	Kenton_Orlin	2017-02-16 18:22:11
2	Andre_Purdy85	2017-04-02 17:11:21
3	Harley_Lind18	2017-02-21 11:11:23
4	Arely_Bogdan63	2016-08-13 03:28:43
5	Aniya_Hackett	2016-12-07 10:46:39
6	Travon_Waters	2017-04-30 13:26:14
7	Kassandra_Homerick	2016-12-12 06:50:08
8	Tabitha_Schamberger11	2016-08-20 02:19:46
9	Gus93	2016-06-24 19:36:31

- Action Output:**

 - 1 12:09:28 create database ig_clone
 - 2 12:09:41 show tables
 - 3 12:11:40 SELECT * FROM users ORDER BY created_at LIMIT 5
 - 4 12:12:09 SELECT * FROM users ORDER BY created_at LIMIT 5
 - 5 12:15:25 table users

Session 2 (Bottom):

 - Query:** `SELECT * FROM users ORDER BY created_at LIMIT 5`
 - Result Grid:**

ID	username	created_at
80	Darby_Herzog	2016-05-06 00:14:21
67	Emilio_Bernier52	2016-05-06 13:04:30
63	Elenor88	2016-05-08 01:30:41
95	Nicole71	2016-05-09 17:30:22
38	Jordyn_Jackson2	2016-05-14 07:56:26

 - Action Output:**

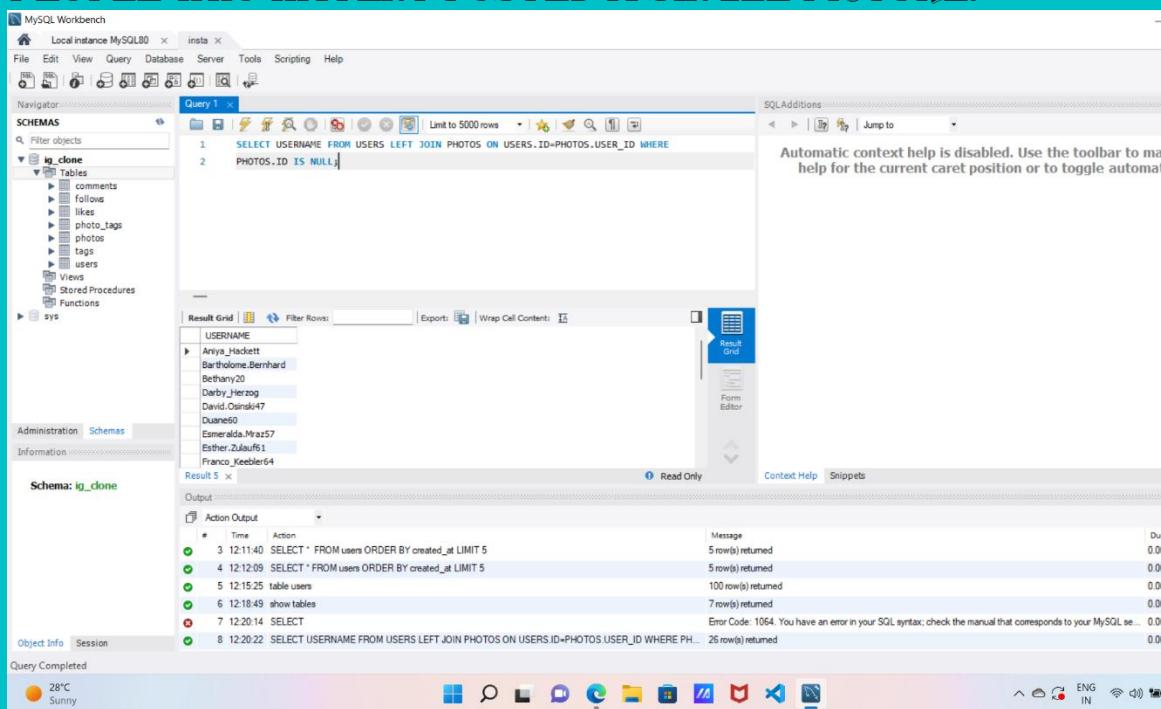
 - 1 12:09:28 create database ig_clone
 - 2 12:09:41 show tables
 - 3 12:11:40 SELECT * FROM users ORDER BY created_at LIMIT 5
 - 4 12:12:09 SELECT * FROM users ORDER BY created_at LIMIT 5

A VERY STRAIGHT FORWARD QUERY WHICH WE EXECUTED BY SELECTING ALL COLUMNS FROM TABLE USERS AND ORDERING IT BY THE DATE THE USERS CREATED THEIR ACCOUNT AND BECAUSE WE NEED 5 OLDEST, WE USED THE LIMIT FUNCTION.

THESE ARE LOYAL TO INSTAGRAM APP AND HENCE WE CAN REWARD THEM FOR THE SAME.

2)

SECOND QUERY WAS ABOUT KNOWING THE INACTIVE PEOPLE WHO HAVENT POSTED A SINGLE PICTURE.



The screenshot shows the MySQL Workbench interface. The left sidebar displays the schema 'ig_clone' with various tables like users, photos, and likes. The main area has a query editor with the following SQL code:

```
1 SELECT USERNAME FROM USERS LEFT JOIN PHOTOS ON USERS.ID=PHOTOS.USER_ID WHERE
2 PHOTOS.ID IS NULL;
```

The results grid shows a list of usernames:

USERNAME
Anya_Hackett
Bartholem_Bernhard
Bethany20
Darby_Herzog
David_Osinski47
Dunne60
Emeralda_Mraz57
Esther_Zuluski1
Franco_Keebler64

The output pane shows the execution history:

Action	Time	Message	Duration
3 12:11:40	SELECT * FROM users ORDER BY created_at LIMIT 5	5 row(s) returned	0.00
4 12:12:00	SELECT * FROM users ORDER BY created_at LIMIT 5	5 row(s) returned	0.00
5 12:15:25	table users	100 row(s) returned	0.00
6 12:18:49	show tables	7 row(s) returned	0.00
7 12:20:14	SELECT	Error Code: 1064. You have an error in your SQL syntax; check the manual that corresponds to your MySQL se...	0.00
8 12:20:22	SELECT USERNAME FROM USERS LEFT JOIN PHOTOS ON USERS.ID=PHOTOS.USER_ID WHERE PH...	26 row(s) returned	0.00

HERE WE NEEDED THE LEFT JOIN AS WE NEEDED THE USERS FROM THE USERS TABLE AND THE MATCHED RECORDS FROM PHOTOS TABLE. USER ID IS IN BOTH THE TABLES HENCE WE USED USER ID TO JOIN THESE TABLES. WHEREVER THE PHOTOS.ID IS NULL I.E NO PICTURE, IT RETURNS THE USERNAME.

WE LEARNED THERE ARE PLENTY OF PEOPLE WHO STILL HAVENT POSTED A SINGLE PICTURE, COULD BE BECAUSE OF THEIR BUSY LIFES, COULD BE BECAUSE THEY MIGHT HAVE SOME PROBLEM WITH THIS APP, SO WE CAN REMIND THEM TO POST.

3)

THE NEXT QUERY WAS TO FIND WHICH USER RECEIVED THE HIGHEST NUMBER OF LIKES.

The screenshot shows the MySQL Workbench interface. In the top-left, the 'Navigator' pane displays the schema 'ig_clone' with tables like 'photos', 'likes', and 'users'. The main area, 'Query 1', contains the following SQL code:

```
1 • SELECT USERNAME,PHOTOS.ID,PHOTOS.IMAGE_URL,COUNT(*) AS TOTAL FROM PHOTOS INNER JOIN LIKES  
2 ON LIKES.PHOTO_ID=PHOTOS.ID INNER JOIN USERS ON PHOTOS.USER_ID=USERS.ID GROUP BY PHOTOS.ID  
3 ORDER BY TOTAL DESC LIMIT 1;
```

The 'Result Grid' shows one row of data:

USERNAME	ID	IMAGE_URL	TOTAL
Zack_Kemmer93	145	https://jarret.name	48

The 'Result 6' pane shows the session history:

Action	Time	Message	Dur
SELECT * FROM users ORDER BY created_at LIMIT 5	4 12:12:09	5 row(s) returned	0.00
table users	5 12:15:25	100 row(s) returned	0.00
show tables	6 12:18:49	7 row(s) returned	0.00
SELECT	7 12:20:14	Error Code: 1064. You have an error in your SQL syntax; check the manual that corresponds to your MySQL se...	0.00
SELECT USERNAME FROM USERS LEFT JOIN PHOTOS ON USERS.ID=PHOTOS.USER_ID WHERE PH...	8 12:20:22	26 row(s) returned	0.00
SELECT USERNAME,PHOTOS.ID,PHOTOS.IMAGE_URL,COUNT(*) AS TOTAL FROM PHOTOS INNER JO...	9 12:25:31	1 row(s) returned	0.01

WE USED THE COUNT FUNCTION TO COUNT THE LIKES AND USED THE RELATED COLUMN I.E PHOTO ID HERE TO JOIN TABLES, LIKES AND USERS TO PHOTOS. GROUP BY WAS USED TO GROUP ROWS THAT HAVE THE SAME VALUES AND ORDER BY THE TOTAL NUMBER OF LIKES. HERE, ZACK RECEIVED THE HIGHEST NUMBER OF LIKES WITH A TOTAL OF 48.

4)

THE 4TH QUERY WAS TO POST THE FIRST FIVE POPULAR HASHTAGS SO THAT WE CAN REACH MOST PEOPLE.

The screenshot shows the MySQL Workbench interface with a query editor and a results grid. The query is:

```
1 • 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 DESC LIMIT 5;
```

The results grid displays the following data:

TAG_NAME	TOTAL
smile	59
beach	42
party	39
fun	38
concert	24

The session history shows the following log entries:

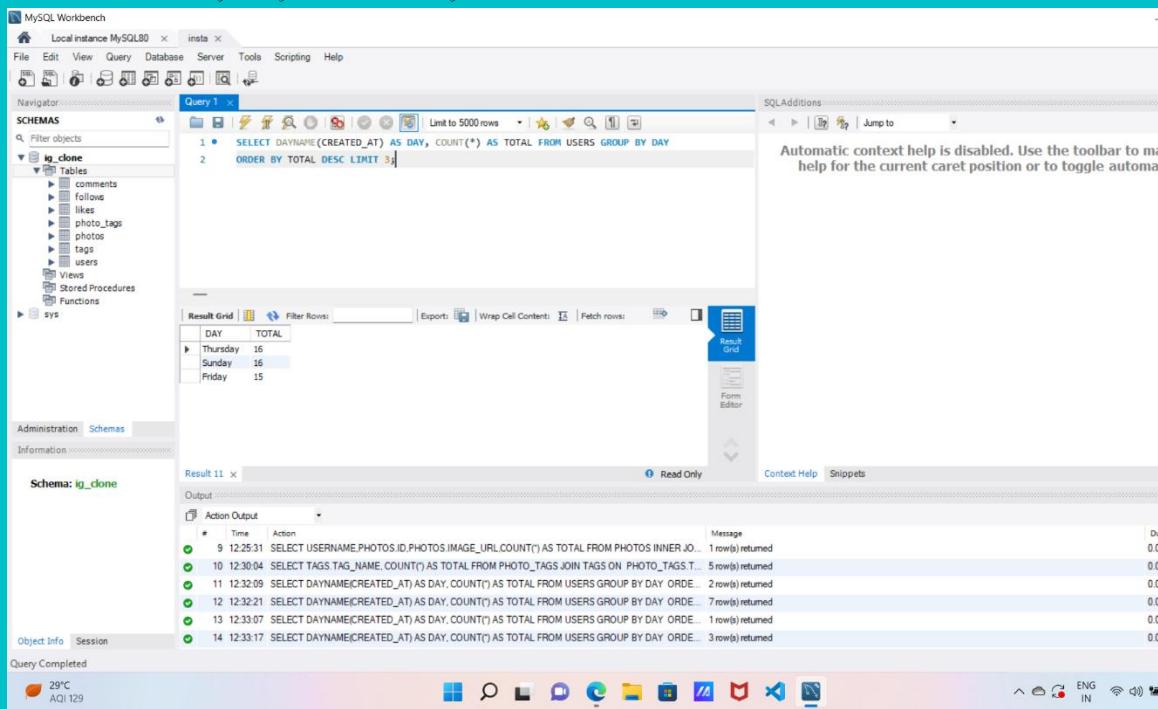
Action	Time	Message	Dur
table users	5 12:15:25	100 row(s) returned	0.00
show tables	6 12:18:49	7 row(s) returned	0.00
SELECT	7 12:20:07	Error Code: 1064. You have an error in your SQL syntax; check the manual that corresponds to your MySQL se... 0.00	
SELECT USERNAME FROM USERS LEFT JOIN PHOTOS ON USERS.ID=PHOTOS.USER_ID WHERE PHOTOS.ID=1	8 12:20:22	26 row(s) returned	0.00
SELECT USERNAME,PHOTOS.ID,PHOTOS.IMAGE_URL,COUNT(*) AS TOTAL FROM PHOTOS INNER JO...	9 12:25:31	1 row(s) returned	0.01
SELECT TAGS.TAG_NAME,COUNT(*) AS TOTAL FROM PHOTO_TAGS JOIN TAGS ON PHOTO_TAGS.T...	10 12:30:04	5 row(s) returned	0.00

SIMILAR TO THE ABOVE QUERY, WE SELECTED THE COUNT OF TAG NAMES, USED TAG ID TO JOIN PHOTOS AND TAGS TABLE , GROUP BY TO GROUP ROWS WITH SAME VALUES I.E WITH SAME TAGS.ID AND ORDER BY TOTAL I.E COUNT OF TAGS NAME IN DESCENDING ORDER SO THE HIGHEST NUMBER COMES ABOVE.

ALL THE HASHTAGS SOUND PRETTY HAPPY, NO WONDER THEY ARE THE MOST USED ONES.

5)

5TH QUERY WAS ABOUT KNOWING WHICH DAY WOULD BE THE BEST TO LAUNCH AD CAMPAIGNS. WHAT DAY DP MOST USERS REGISTER ON THIS APP?



The screenshot shows the MySQL Workbench interface. In the top-left, the Navigator pane displays the schema 'ig_clone' with tables like 'comments', 'follow', 'likes', 'photo_tags', 'photos', 'tags', and 'users'. The main area contains two tabs: 'Query 1' and 'Result 11'. The 'Query 1' tab shows the following SQL code:

```
1 SELECT DAYNAME(CREATED_AT) AS DAY, COUNT(*) AS TOTAL FROM USERS GROUP BY DAY
2 ORDER BY TOTAL DESC LIMIT 3;
```

The 'Result 11' tab shows the output of the query:

DAY	TOTAL
Thursday	16
Sunday	16
Friday	15

Below the results, the 'Output' tab displays the execution history of the session:

Action	Time	Action	Message	Duration
9	12:25:31	SELECT USERNAME.PHOTOS.ID PHOTOS.IMAGE_URL,COUNT(*) AS TOTAL FROM PHOTOS INNER JO...	1 row(s) returned	0.01
10	12:30:04	SELECT TAGS.TAG_NAME,COUNT(*) AS TOTAL FROM PHOTO_TAGS JOIN TAGS ON_PHOTO_TAGS.T...	5 row(s) returned	0.00
11	12:32:09	SELECT DAYNAME(CREATED_AT) AS DAY,COUNT(*) AS TOTAL FROM USERS GROUP BY DAY ORDE...	2 row(s) returned	0.00
12	12:32:21	SELECT DAYNAME(CREATED_AT) AS DAY,COUNT(*) AS TOTAL FROM USERS GROUP BY DAY ORDE...	7 row(s) returned	0.00
13	12:33:07	SELECT DAYNAME(CREATED_AT) AS DAY,COUNT(*) AS TOTAL FROM USERS GROUP BY DAY ORDE...	1 row(s) returned	0.00
14	12:33:17	SELECT DAYNAME(CREATED_AT) AS DAY,COUNT(*) AS TOTAL FROM USERS GROUP BY DAY ORDE...	3 row(s) returned	0.00

WE SELECT THE DAYNAME OF WHEN THE USER CREATED HIS ACCOUNT(USING CREATED_AT), COUNT FUNCTION TO COUNT THE DAYS, GROUP BY TO GROUP ROWS WITH SAME VALUES I.E GROUPING DAYS AND ORDER BY TOTAL IN DESCENDING. WE PRINTED TOP THREE ROWS AND GOT TO KNOW BOTH THURSDAY AND SUNDAY WE HAVE HIGHEST NUMBER OF REGISTRATIONS , 16.

INVESTOR METRICS

1)

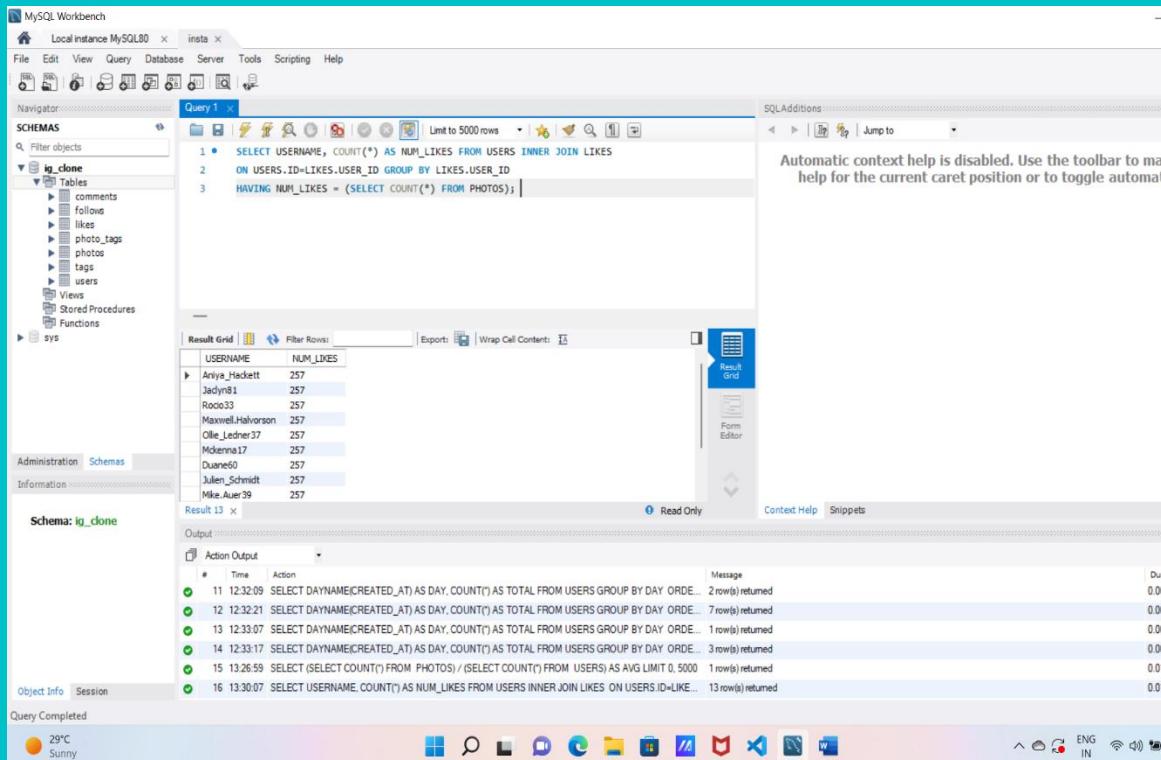
OUR FIRST QUERY HERE WAS TO FIND THE AVERAGE NUMBER OF PHOTOS PER USER

The screenshot shows the MySQL Workbench interface. In the top navigation bar, 'File', 'Edit', 'View', 'Query', 'Database', 'Server', 'Tools', 'Scripting', and 'Help' are visible. Below the toolbar, there's a 'Navigator' pane showing 'SCHEMAS' and 'Tables' for the 'ig_clone' schema, which contains tables like 'comments', 'follows', 'likes', 'photo_tags', 'photos', 'tags', and 'users'. The main area is titled 'Query 1' and contains the SQL query: `SELECT (SELECT COUNT(*) FROM PHOTOS) / (SELECT COUNT(*) FROM USERS) AS AVG;`. The result grid shows a single row with 'AVG' and the value '2.5700'. Below the query results, the 'Result 12' tab shows the history of actions taken in the session. At the bottom, the status bar displays '29°C Sunny' and system icons.

WE DO THIS BY USING THE FORMULA FOR AVERAGE I.E
COUNT OF PHOTOS DIVIDED BY COUNT OF USERS.
THE AVERAGE COMES OUT TO BE 25700

2)

THE LAST QUERY WAS TO FIND THE BOT ACCOUNTS, WE KNOW THIS BY FINDING WHO ALL HAS LIKED ALL PHOTOS SINCE NORMAL USERS CANNOT DO THIS.



The screenshot shows the MySQL Workbench interface with a query editor and results grid. The query is:

```
1 • SELECT USERNAME, COUNT(*) AS NUM_LIKES FROM USERS INNER JOIN LIKES  
2   ON USERS.ID=LIKES.USER_ID GROUP BY LIKES.USER_ID  
3   HAVING NUM_LIKES = (SELECT COUNT(*) FROM PHOTOS);
```

The results grid shows the following data:

USERNAME	NUM_LIKES
Anya_Hackett	257
Jaclyn81	257
Roco33	257
MarkwellHalvorson	257
Ollie_Ledner37	257
Mikenna17	257
Duanee60	257
Julien_Schmidt	257
Mike_Auer39	257

The output pane shows the execution history:

Action	Time	Action	Message	Duration
11	12:32:06	SELECT DAYNAME(CREATED_AT) AS DAY, COUNT(*) AS TOTAL FROM USERS GROUP BY DAY ORDER BY DAY	2 row(s) returned	0.00
12	12:32:24	SELECT DAYNAME(CREATED_AT) AS DAY, COUNT(*) AS TOTAL FROM USERS GROUP BY DAY ORDER BY DAY	7 row(s) returned	0.00
13	12:33:07	SELECT DAYNAME(CREATED_AT) AS DAY, COUNT(*) AS TOTAL FROM USERS GROUP BY DAY ORDER BY DAY	1 row(s) returned	0.00
14	12:33:17	SELECT DAYNAME(CREATED_AT) AS DAY, COUNT(*) AS TOTAL FROM USERS GROUP BY DAY ORDER BY DAY	3 row(s) returned	0.00
15	13:26:59	SELECT (SELECT COUNT(*) FROM PHOTOS) / (SELECT COUNT(*) FROM USERS) AS AVG LIMIT 0, 5000	1 row(s) returned	0.01
16	13:30:07	SELECT USERNAME, COUNT(*) AS NUM_LIKES FROM USERS INNER JOIN LIKES ON USERS.ID=LIKES.USER_ID GROUP BY LIKES.USER_ID HAVING NUM_LIKES = (SELECT COUNT(*) FROM PHOTOS)	13 row(s) returned	0.01

WE USE COUNT FUNCTION TO COUNT THE NUMBER OF LIKES FROM USERS WHICH IS EQUAL TO THE TOTAL NUMBER OF PHOTOS. USING USER_ID TO JOIN TABLES USERS AND LIKES, WE USE A SUB QUERY TO FIND OUT THE TOTAL NUMBER OF PHOTOS.

AS WE CAN SEE THERE ARE PLENTY OF BOT ACCOUNTS WHICH NEEDS TO BE ERADICATED.

OPERATION ANALYTICS AND INVESTIGATING METRIC SPIKE

PROJECT DESCRIPTION

OPERATION ANALYTICS IS THE ANALYSIS DONE FOR THE COMPLETE END TO END OPERATIONS OF A COMPANY. WITH THE HELP OF THIS, THE COMPANY THEN FINDS THE AREAS ON WHICH IT MUST IMPROVE UPON, WORKING CLOSELY WITH THE OPS TEAM, SUPPORT TEAM, MARKETING TEAM AND HELP THEM DERIVE INSIGHTS OUT OF THE DATA THEY COLLECT. BEING ONE OF THE MOST IMPORTANT PARTS OF A COMPANY, THIS KIND OF ANALYSIS IS FURTHER USED TO PREDICT THE OVERALL GROWTH OR DECLINE OF A COMPANY'S FORTUNE. IT ALSO USED FOR BETTER AUTOMATION, BETTER UNDERSTANDING BETWEEN CROSS-FUNCTIONAL TEAMS, AND MORE EFFECTIVE WORKFLOWS.

INVESTIGATING METRIC SPIKE IS ALSO AN IMPORTANT PART OF OPERATION ANALYTICS AS BEING A DATA ANALYST YOU MUST BE ABLE TO UNDERSTAND OR MAKE OTHER TEAMS UNDERSTAND QUESTIONS LIKE- WHY IS THERE A DIP IN DAILY ENGAGEMENT? WHY HAVE SALES TAKEN A DIP? ETC. QUESTIONS LIKE THESE MUST BE ANSWERED DAILY AND FOR THAT ITS VERY IMPORTANT TO INVESTIGATE METRIC SPIKE.

APPROACH

I GATHERED INFORMATION FROM THE DESCRIPTION. TOOK A LOOK AT THE DATASET TO UNDERSTAND ABOUT THE TABLES AND THEN STARTED WORKING ON THIS PROJECT. THEN I USED MYSQL WORKBENCH TO CREATE TABLES AND FINALLY EXECUTED MY QUERIES.

A FEW OF THEM HAD ERRORS BUT I REVISED MY QUERY AND EVENTUALLY GOT THEM RIGHT.

INSIGHTS

THE KNOWLEDGE THAT I GAINED FROM THIS PROJECT IS HOW TO CREATE THE DATABASE ,HOW TO USE THE DATABASE , HOW CAN WE PERFORM SQL OPERATION BASED ON SCENARIOS AND HOW WE CAN USE THIS DATA TO IMPROVE THE EXPERIENCE ALTOGETHER WHILE HELPING THE BUSINESS GROW.

JOB DATA

A NUMBER OF JOBS REVIEWED: AMOUNT OF JOBS REVIEWED OVER TIME.

YOUR TASK: CALCULATE THE NUMBER OF JOBS REVIEWED PER HOUR PER DAY FOR NOVEMBER 2020?

DS	THROUGHPUT
30-11-2020	144
29-11-2020	180
28-11-2020	218.18
27-11-2020	34.62
25-11-2020	80

B.THROUGHPUT: IT IS THE NO. OF EVENTS HAPPENING PER SECOND.

YOUR TASK: LET'S SAY THE ABOVE METRIC IS CALLED THROUGHPUT. CALCULATE 7 DAY ROLLING AVERAGE OF THROUGHPUT? FOR THROUGHPUT, DO YOU PREFER DAILY METRIC OR 7-DAY ROLLING AND WHY?

DS	THROUGHPUT_FOR_7_DAYS
25-11-2020	0.02
27-11-2020	0.01
28-11-2020	0.02
29-11-2020	0.02
30-11-2020	0.03

C.PERCENTAGE SHARE OF EACH LANGUAGE: SHARE OF EACH LANGUAGE FOR DIFFERENT CONTENTS.

YOUR TASK: CALCULATE THE PERCENTAGE SHARE OF EACH LANGUAGE IN THE LAST 30 DAYS?

LANGUAGE	PERCENTAGE_OF_JOBS
PERSIAN	200
ITALIAN	100
FRENCH	100
HINDI	100
ARABIC	100
PERSIAN	100
ITALIAN	50
FRENCH	50
HINDI	50
ARABIC	50

D.DUPLICATE ROWS: ROWS THAT HAVE THE SAME VALUE PRESENT IN THEM.

YOUR TASK: LET'S SAY YOU SEE SOME DUPLICATE ROWS IN THE DATA.

HOW WILL YOU DISPLAY DUPLICATES FROM THE TABLE?

NO DUPLICATE DATA FOUND IN THE TABLE

Ds	job_id	actor_id	Event	language	time_spent	org

INVESTIGATING METRIC SPIKE

A. USER ENGAGEMENT: TO MEASURE THE ACTIVENESS OF A USER.
MEASURING IF THE USER FINDS QUALITY IN A PRODUCT/SERVICE

YOUR TASK: CALCULATE THE WEEKLY USER ENGAGEMENT?

DATE_TRUNC	WEEKLY_ACTIVE_USERS
28-04-2014 00:00	701
05-05-2014 00:00	1054
12-05-2014 00:00	1094
19-05-2014 00:00	1147
26-05-2014 00:00	1113
02-06-2014 00:00	1173
09-06-2014 00:00	1219
16-06-2014 00:00	1262
23-06-2014 00:00	1249
30-06-2014 00:00	1271
07-07-2014 00:00	1355
14-07-2014 00:00	1345
21-07-2014 00:00	1363
28-07-2014 00:00	1442
04-08-2014 00:00	1266
11-08-2014 00:00	1215

18-08-2014 00:00 1203

25-08-2014 00:00 1194

B.USER GROWTH: AMOUNT OF USERS GROWING OVER TIME FOR A PRODUCT.

YOUR TASK: CALCULATE THE USER GROWTH FOR PRODUCT?

DAY	ALL_USE_RS	ACTIVATED_USERS
01-01-2013 00:00	13	7
02-01-2013 00:00	11	7
03-01-2013 00:00	14	6
04-01-2013 00:00	11	1
05-01-2013 00:00	3	2
06-01-2013 00:00	4	3
07-01-2013 00:00	13	4
08-01-2013 00:00	13	2
09-01-2013 00:00	11	6
10-01-2013 00:00	12	6
11-01-2013 00:00	11	6
12-01-2013 00:00	4	3
13-01-2013 00:00	3	2
14-01-2013 00:00	13	8
15-01-2013 00:00	15	11
16-01-2013 00:00	14	7
17-01-2013 00:00	13	9
18-01-2013 00:00	14	10
19-01-2013 00:00	4	1
20-01-2013 00:00	4	1
21-01-2013 00:00	13	7
22-01-2013 00:00	13	5
23-01-2013 00:00	14	7
24-01-2013 00:00	14	5
25-01-2013 00:00	16	8
26-01-2013 00:00	3	3
27-01-2013 00:00	4	1
28-01-2013 00:00	14	7

29-01-2013 00:00	14	3
30-01-2013 00:00	14	6

C.WEEKLY RETENTION: USERS GETTING RETAINED WEEKLY AFTER SIGNING-UP FOR A PRODUCT.

YOUR TASK: CALCULATE THE WEEKLY RETENTION OF USERS-SIGN UP COHORT?

WEEK	AVERAGE AGE DURING WEEK	10+ WEEKS		9 WEEKS		8 WEEKS		7 WEEKS		6 WEEKS		5 WEEKS		4 WEEKS		3 WEEKS		
		S	S	S	S	K	S	K	S	K	S	K	S	K	S	K	S	K
28-04-2014 00:00	124.007238																	
	9	701	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
05-05-2014 00:00	124.381690																	
	8	1054	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
12-05-2014 00:00	131.938644																	
	2	1094	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
19-05-2014 00:00	132.326628																	
	4	1147	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
26-05-2014 00:00	132.345363																	
	4	1113	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
02-06-2014 00:00	131.831109																	
	1	1173	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
09-06-2014 00:00	131.042582																	
	4	1219	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
16-06-2014 00:00	136.480565																	
	4	1255	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
23-06-2014 00:00	136.278905																	
	6	1034	210	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
30-06-2014 00:00	136.419297																	
	5	917	151	199	0	0	0	0	0	0	0	0	0	0	0	0	0	0
07-07-2014 00:00	135.888750																	
	5	899	100	130	223	0	0	0	0	0	0	0	0	0	0	0	0	0
14-07-2014 00:00	143.448815																	
	7	832	62	82	152	215	0	0	0	0	0	0	0	0	0	0	0	0
21-07-2014 00:00	141.70278	791	44	60	95	144	228	0	0	0	0	0	0	0	0	0	0	0
28-07-2014 00:00	144.078660																	
	4	805	30	43	83	91	155	234	0	0	0	0	0	0	0	0	0	0

04-08-2014											
00:00	140.732238	678	24	34	52	52	82	154	189	0	0
11-08-2014	125.994310										
00:00	1	562	19	33	39	33	59	94	126	0	0
18-08-2014	128.021718										
00:00	1	522	15	26	26	19	40	64	69	59	0
25-08-2014	128.269810										
00:00	4	474	15	14	23	20	31	47	48	73	20

D.WEEKLY ENGAGEMENT: TO MEASURE THE ACTIVENESS OF A USER.

MEASURING IF THE USER FINDS QUALITY IN PRODUCT/SERVICE_WEEKLY.

YOUR TASK: CALCULATE THE WEEKLY ENGAGEMENT PER DEVICE?

WEEK	WEEKLY_ACTIVE_USER_S	COMPUTER	PHONE	TABLET
28-04-2014 00:00	701	347	211	103
05-05-2014 00:00	1054	590	335	176
12-05-2014 00:00	1094	632	352	191
19-05-2014 00:00	1147	647	385	181
26-05-2014 00:00	1113	605	379	176
02-06-2014 00:00	1173	663	388	197
09-06-2014 00:00	1219	671	397	195
16-06-2014 00:00	1262	685	400	227
23-06-2014 00:00	1249	706	395	210
30-06-2014 00:00	1271	676	424	218
07-07-2014 00:00	1355	735	434	227
14-07-2014 00:00	1345	762	421	218
21-07-2014 00:00	1363	767	436	218
28-07-2014 00:00	1442	798	414	241
04-08-2014 00:00	1266	770	354	166
11-08-2014 00:00	1215	722	326	153
18-08-2014 00:00	1203	754	311	145
25-08-2014 00:00	1194	725	312	150

D.EMAIL ENGAGEMENT: USERS ENGAGING WITH THE EMAIL SERVICE.

YOUR TASK: CALCULATE THE EMAIL ENGAGEMENT METRICS?

WEEK	WEEKLY_EMAILS	REENGAGEMENT_EMAIL	EMAIL_OPENS	EMAIL_CLICKS
14-07-2014 00:00	3499	226	1260	607
28-07-2014 00:00	3706	230	1386	633
19-05-2014 00:00	2733	179	995	498
05-05-2014 00:00	2602	164	919	434
11-08-2014 00:00	3897	224	1357	430
12-05-2014 00:00	2665	175	971	479
25-08-2014 00:00	4111	263	1533	493
28-04-2014 00:00	908	98	332	187
26-05-2014 00:00	2822	179	1026	453
07-07-2014 00:00	3399	214	1230	622
09-06-2014 00:00	3003	190	1070	533
16-06-2014 00:00	3105	234	1161	563
30-06-2014 00:00	3302	222	1168	559
21-07-2014 00:00	3592	206	1211	584
02-06-2014 00:00	2911	199	993	492
18-08-2014 00:00	4012	257	1421	487
23-06-2014 00:00	3207	187	1090	524
04-08-2014 00:00	3793	206	1336	432

HIRING PROCESS ANALYTICS

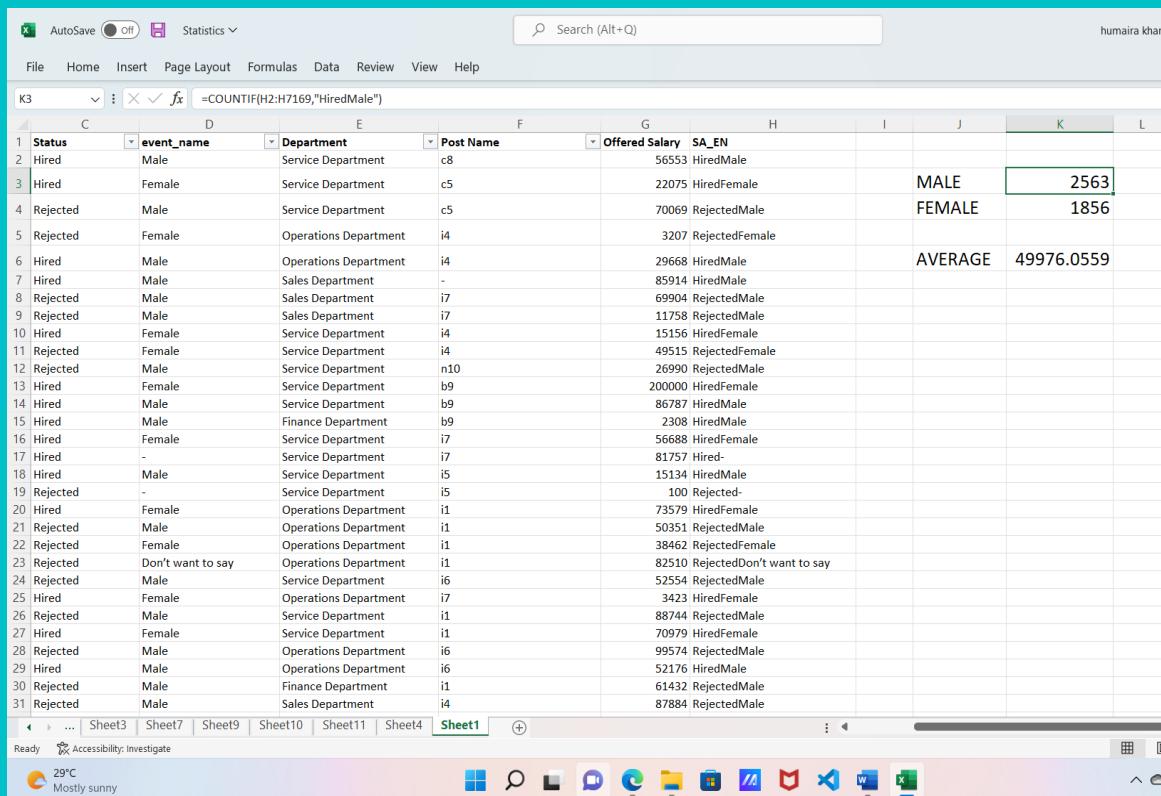
PROJECT DESCRIPTION: THIS PROJECT IS TO KNOW ABOUT THE MAJOR UNDERLYING TRENDS ABOUT THE HIRING PROCESS. TRENDS SUCH AS- NUMBER OF REJECTIONS, NUMBER OF INTERVIEWS, TYPES OF JOBS, VACANCIES ETC. ARE IMPORTANT FOR A COMPANY TO ANALYZE BEFORE HIRING FRESHERS OR ANY OTHER INDIVIDUAL. BEING A DATA ANALYST, OUR JOB IS TO GO THROUGH THESE TRENDS AND DRAW INSIGHTS OUT OF IT FOR HIRING DEPARTMENT TO WORK UPON.

APPROACH AND INSIGHTS

I FIRST DOWNLOADED THE DATASET PROVIDED AND OPENED IT IN EXCEL. TRIED TO UNDERSTAND WHAT THAT DATASET CONTAINS, AND THEN STARTED WITH THE QUERIES USING ONLY EXCEL SHEET.

THE FIRST TASK WAS TO FIND THE NUMBER OF MALE AND FEMALE WORKERS HIRED. SO MY FIRST INSTINCT WAS TO JOIN STATUS COLUMN AND EVENT NAME COLUMN SO THAT THE NEW COLUMN CREATED HAS DETAILS OF BOTH WHO HAS BEEN HIRED ALONG WITH THE GENDER.

I USED THE COUNTIF FUNCTION,
`=COUNTIF(G2:G7169,"HIREDMALE")`
`=COUNTIF(G2:G7169,"HIREDFEMALE")`
BELOW IS THE ANSWER.



C	D	E	F	G	H	I	J	K	L
Status	event_name	Department	Post Name	Offered Salary	SA_EN				
2 Hired	Male	Service Department	c8	56553	HiredMale				
3 Hired	Female	Service Department	c5	22075	HiredFemale	MALE	2563		
4 Rejected	Male	Service Department	c5	70069	RejectedMale	FEMALE	1856		
5 Rejected	Female	Operations Department	i4	3207	RejectedFemale				
6 Hired	Male	Operations Department	i4	29668	HiredMale	AVERAGE	49976.0559		
7 Hired	Male	Sales Department	-	85914	HiredMale				
8 Rejected	Male	Sales Department	i7	69904	RejectedMale				
9 Rejected	Male	Sales Department	i7	11758	RejectedMale				
10 Hired	Female	Service Department	i4	15156	HiredFemale				
11 Rejected	Female	Service Department	i4	49515	RejectedFemale				
12 Rejected	Male	Service Department	n10	26990	RejectedMale				
13 Hired	Female	Service Department	b9	200000	HiredFemale				
14 Hired	Male	Service Department	b9	86787	HiredMale				
15 Hired	Male	Finance Department	b9	2308	HiredMale				
16 Hired	Female	Service Department	i7	56688	HiredFemale				
17 Hired	-	Service Department	i7	81757	Hired-				
18 Hired	Male	Service Department	i5	15134	HiredMale				
19 Rejected	-	Service Department	i5	100	Rejected-				
20 Hired	Female	Operations Department	i1	73579	HiredFemale				
21 Rejected	Male	Operations Department	i1	50351	RejectedMale				
22 Rejected	Female	Operations Department	i1	38462	RejectedFemale				
23 Rejected	Don't want to say	Operations Department	i1	82510	RejectedDon't want to say				
24 Rejected	Male	Service Department	i6	52554	RejectedMale				
25 Hired	Female	Operations Department	i7	3423	HiredFemale				
26 Rejected	Male	Service Department	i1	88744	RejectedMale				
27 Hired	Female	Service Department	i1	70979	HiredFemale				
28 Rejected	Male	Operations Department	i6	99574	RejectedMale				
29 Hired	Male	Operations Department	i6	52176	HiredMale				
30 Rejected	Male	Finance Department	i1	61432	RejectedMale				
31 Rejected	Male	Sales Department	i4	87884	RejectedMale				

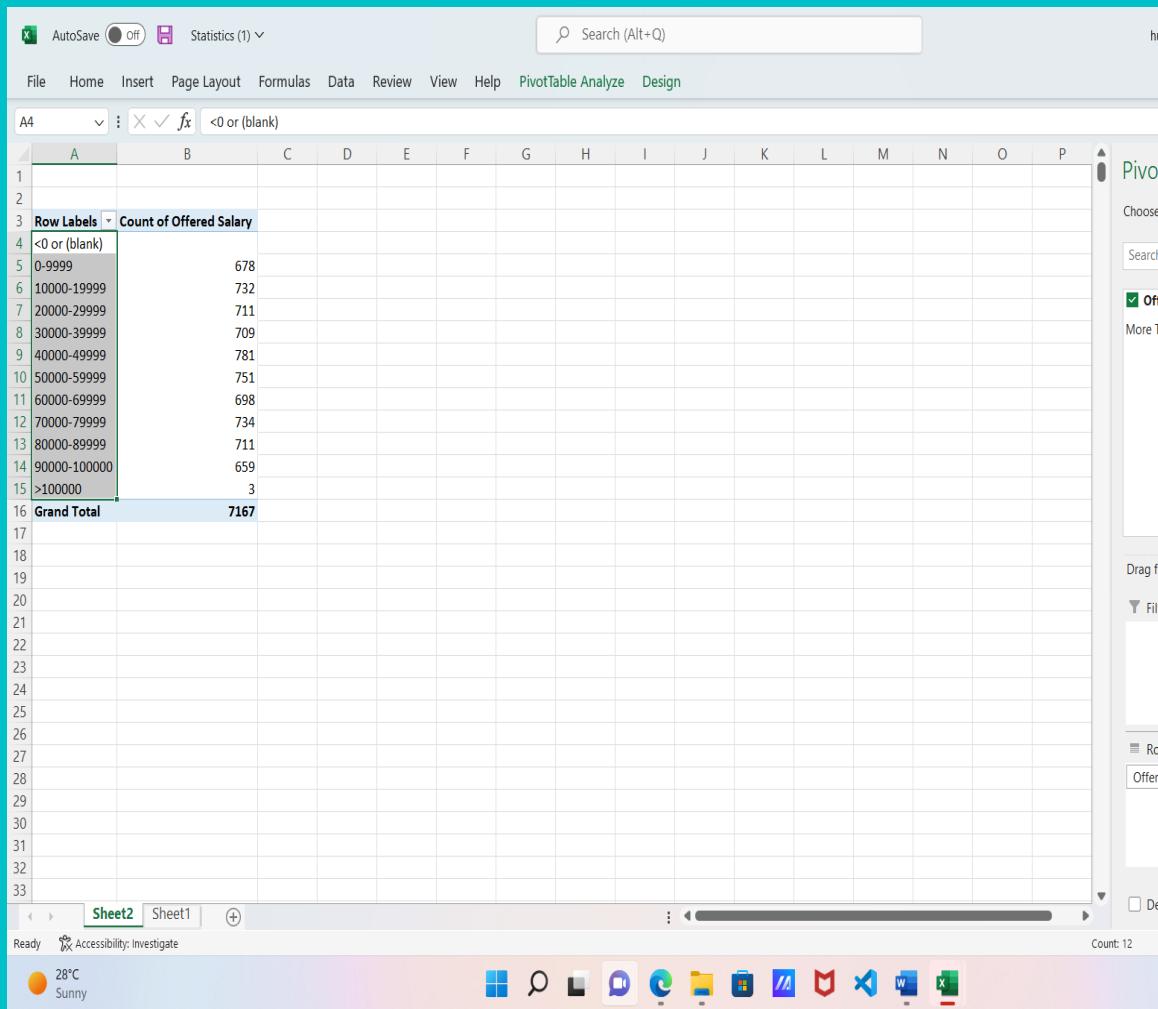
2563 MALES AND 1856 FEMALES WERE HIRED.

**SECOND QUERY WAS TO FIND THE AVERAGE SALARY.
AS YOU CAN SEE IN THE PICTURE, AVERAGE COMES OUT
TO BE 49976.0559 USING FORMULA
=AVERAGE(G2:G7169)**

The screenshot shows a Microsoft Excel spreadsheet titled "Statistics". The formula bar at the top displays the formula `=COUNTIF(H2:H7169,"HiredMale")`. The main area contains a data table with columns: Status, event_name, Department, Post Name, Offered Salary, and SA_EN. The data includes various rows for Hired and Rejected individuals across different departments like Service, Sales, and Finance. In the bottom right corner of the data area, the average salary is calculated as `AVERAGE G2:G7169`, resulting in `49976.0559`.

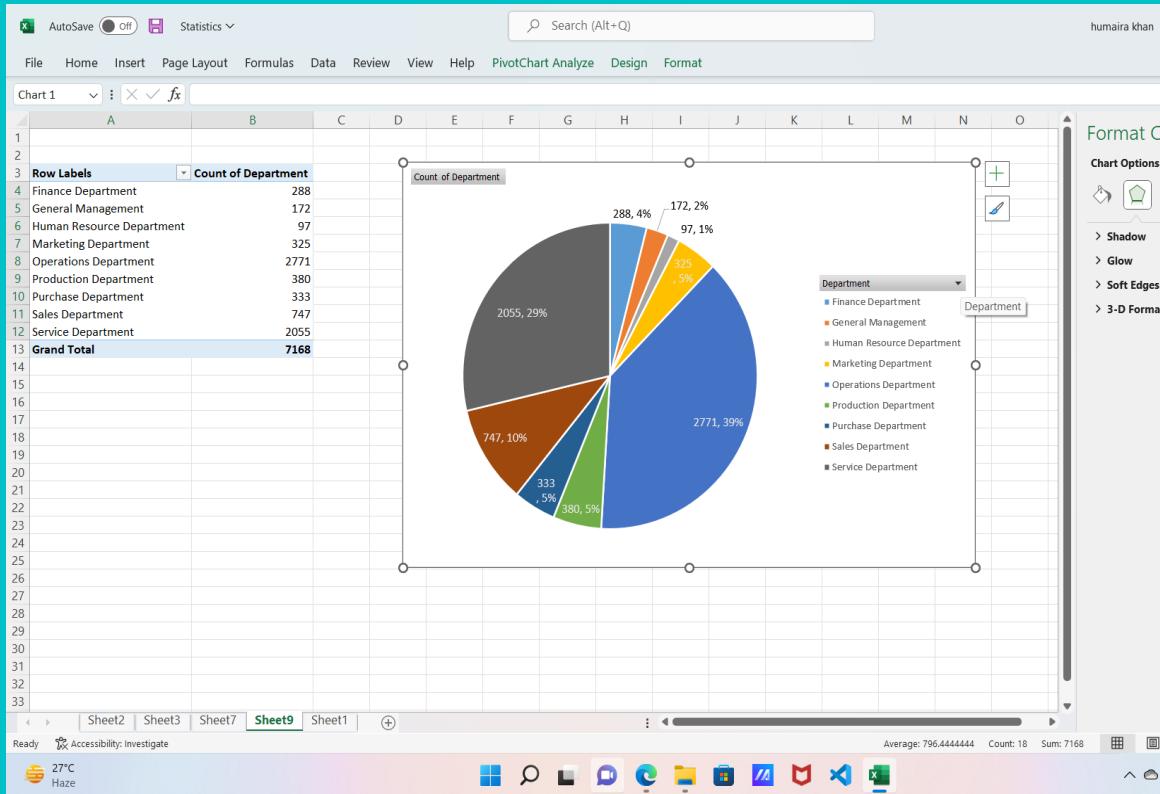
C	D	E	F	G	H	I	J	K	L
1 Status	event_name	Department	Post Name	Offered Salary	SA_EN				
2 Hired	Male	Service Department	c8	56553	HiredMale				
3 Hired	Female	Service Department	c5	22075	HiredFemale	MALE	2563		
4 Rejected	Male	Service Department	c5	70069	RejectedMale	FEMALE	1856		
5 Rejected	Female	Operations Department	i4	3207	RejectedFemale				
6 Hired	Male	Operations Department	i4	29668	HiredMale	AVERAGE	49976.0559		
7 Hired	Male	Sales Department	-	85914	HiredMale				
8 Rejected	Male	Sales Department	i7	69904	RejectedMale				
9 Rejected	Male	Sales Department	i7	11758	RejectedMale				
10 Hired	Female	Service Department	i4	15156	HiredFemale				
11 Rejected	Female	Service Department	i4	49515	RejectedFemale				
12 Rejected	Male	Service Department	n10	26990	RejectedMale				
13 Hired	Female	Service Department	b9	200000	HiredFemale				
14 Hired	Male	Service Department	b9	86787	HiredMale				
15 Hired	Male	Finance Department	b9	2308	HiredMale				
16 Hired	Female	Service Department	i7	56688	HiredFemale				
17 Hired	-	Service Department	i7	81757	Hired-				
18 Hired	Male	Service Department	i5	15134	HiredMale				
19 Rejected	-	Service Department	i5	100	Rejected-				
20 Hired	Female	Operations Department	i1	73579	HiredFemale				
21 Rejected	Male	Operations Department	i1	50351	RejectedMale				
22 Rejected	Female	Operations Department	i1	38462	RejectedFemale				
23 Rejected	Don't want to say	Operations Department	i1	82510	RejectedDon't want to say				
24 Rejected	Male	Service Department	i6	52554	RejectedMale				
25 Hired	Female	Operations Department	i7	3423	HiredFemale				
26 Rejected	Male	Service Department	i1	88744	RejectedMale				
27 Hired	Female	Service Department	i1	70979	HiredFemale				
28 Rejected	Male	Operations Department	i6	99574	RejectedMale				
29 Hired	Male	Operations Department	i6	52176	HiredMale				
30 Rejected	Male	Finance Department	i1	61432	RejectedMale				
31 Rejected	Male	Sales Department	i4	87884	RejectedMale				

THIRD QUERY WAS TO DRAW CLASS INTERVALS FOR SALARY. I USED PIVOT TABLE AND THE GROUPING FUNCTION TO DO THE SAME.

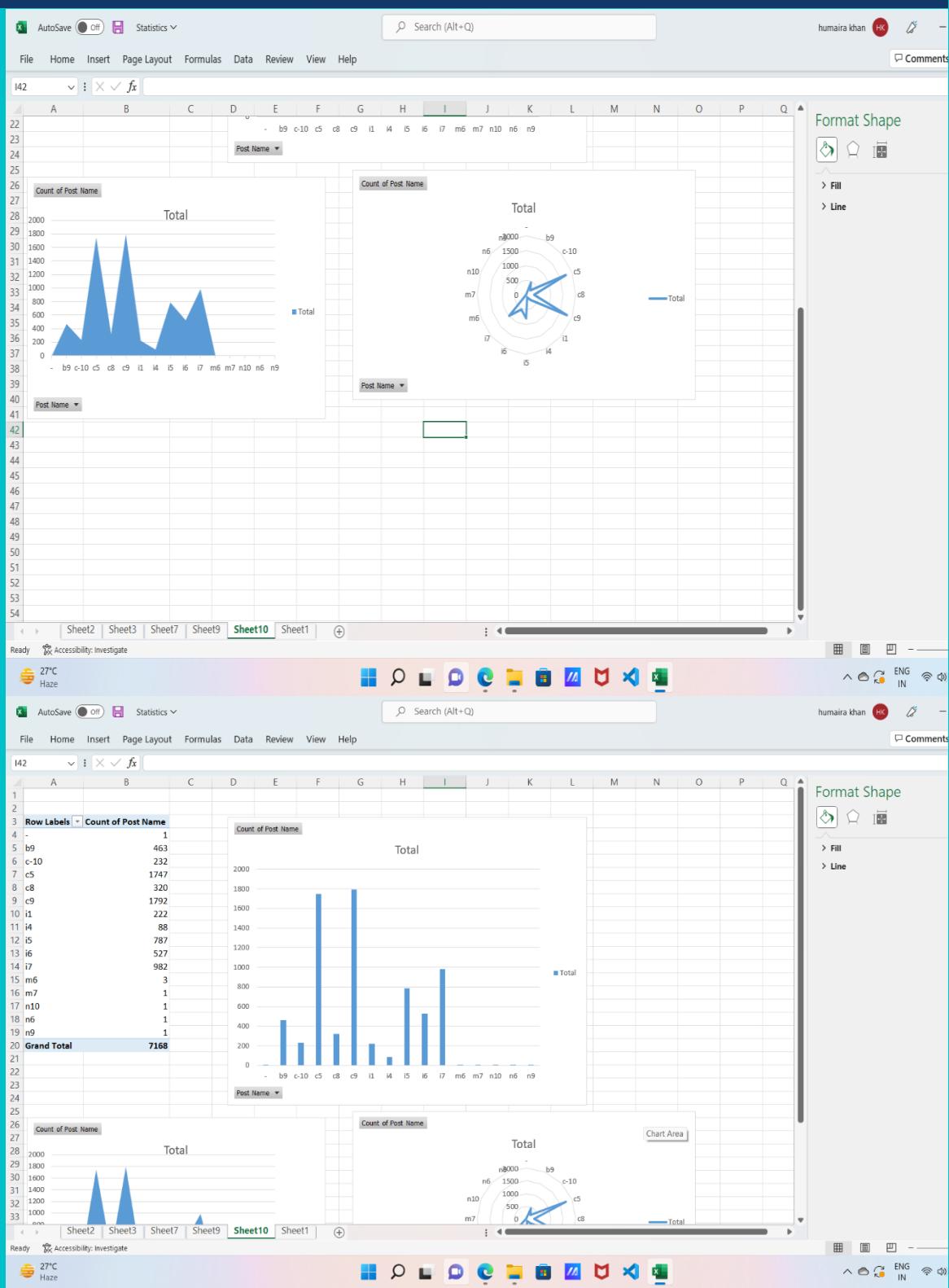


SET THE RANGE TO 10000 RANGING FROM 0 TO 100000.

4TH QUERY WAS TO DRAW PIE CHART TO SHOW PROPORTION OF PEOPLE WORKING DIFFERENT DEPARTMENT.



**AGAIN I USED THE PIVOT TABLE ON THE DEPARTMENT COLUMN FIRST AND THEN CHOSE A PIE CHART. I ADDED LABELS TO THE CHART SO IT BECOMES EASIER TO UNDERSTAND ALONG WITH DATA LABELS.
LAST QUERY WAS TO USE DIFFERENT CHARTS.**



I USED IT ON POST NAMES AND THEIR TOTAL COUNT AND CREATED 3 GRAPHS NAMELY BAR CHART, AREA CHART AND RADAR CHART.

MOVIE IDMB ANALYSIS

DESCRIPTION:

THE PROJECT IS ABOUT FINDING OUT THE VARIOUS INSIGHTS IN IMDB_MOVIES DATASET . WE ANALYZE THIS DATA AND SOME FOLLOWING QUESTIONS:

- 1.CLEAN THE DATA
- 2.FIND THE MOVIES WITH THE HIGHEST PROFIT?
- 3.FIND IMDB TOP 250
- 4.FIND TOP 10 DIRECTORS
- 5.FIND POPULAR GENRES
- 6.FIND THE CRITIC-FAVORITE AND AUDIENCE-FAVORITE ACTORS.

QUERIES

THE FIRST TASK WAS TO CLEAN THE DATASET.
I STARTED THIS WITH REMOVING DUPLICATE DATA.

color	director	n_num	critic	duration	director_f	actor_3	f_actor_2	n_actor_1	f_gross	genres	actor_1	n_movie_titl	num_vote	cast_total	actor_3_n	facenumb	plot	keyw	movie_imdbnum	user_language	count
Color	James Cam	723	178	0	855	Joel David	1000	7.61E+08	Action Ad CCH Pounder Avatar	886204	4834	Wes Studi	0	avatar fut	http://ww	3054	English	US			
Color	Gore Verbi	302	169	563	1000	Orlando Bl	4000	3.09E+08	Action Ad Johnny De Pirates of t	471220	48350	Jack Daven	0	goddess h	http://ww	1238	English	US			
Color	Sam Mend	602	148	0	161	Rory Kinne	11000	2E+08	Action Ad Christoph Spectre	275868	11700	Stephanie	1	bomb esp	http://ww	994	English	UK			
Color	Christophe	813	164	22000	23000	Christian B	27000	4.48E+08	Action Thi Tom Hardy The Dark K	1144337	106759	Joseph Go	0	deception	http://ww	2701	English	US			
Doug Walker				131		Rob Walk	131		Document Doug Wall Star Wars:	8	143		0								
Color	Andrew St.	462	132	475	530	Samantha	640	73058679	Action Ad Daryl Sabo John Carte	212204	1873	Polly Walk	1	alien ame	http://ww	738	English	US			
Color	Sam Raimi	392	156	0	4000	James Frar	24000	3.37E+08	Action Ad J.K. Simmc Spider-Ma	383056	46055	Kirsten Du	0	sandman t	http://ww	1902	English	US			
Color	Nathan Gr	324	100	15	284	Donna Mu	790	2.01E+08	Adventure Brad Garre Tangled	294810	2036	M.C. Gain	1	17th centu	http://ww	387	English	US			
Color	Joss Whed	635	141	0	19000	Robert Do	26000	4.59E+08	Action Ad Chris Hem Avengers:	462669	92000	Scarlett Jo	4	artificial in	http://ww	1117	English	US			
Color	David Yate	375	153	282	10000	Daniel Rad	25000	3.02E+08	Adventure Alan Rickr Harry Pott	321795	58753	Rupert Gri	3	blood boc	http://ww	973	English	UK			
Color	Zack Snyde	673	183	0	2000	Lauren Col	15000	3.3E+08	Action Ad Henry Cav Batman v	371639	24450	Alan D. Pu	0	based on c	http://ww	3018	English	US			
Color	Bryan Sing	434	169	0	903	Marlon Br	100000	3E+08	Action Ad Louis Cse Superman	340205	30001	Frank Lang	0	central epi	http://ww	2367	English	US			
Color	Marc Forst	403	106	395	393	Mathie	Microsoft Excel						X	her http://ww		1243	English	UK			
Color	Gore Verbi	313	151	563	1000	Orland							ice http://ww		1832	English	US				
Color	Gore Verbi	450	150	563	1000	Ruth W							put http://ww		711	English	US				
Color	Zack Snyde	733	143	0	748	Christo							an c http://ww		2536	English	US				
Color	Andrew Ac	258	150	80	201	Pierfr							br http://ww		438	English	US				
Color	Joss Whed	703	173	0	19000	Robert Do	20000	0.25E+08	Action Ad Chris Hem The Aven	555415	67097	Scarlett Jo	3	alien ivas	http://ww	1722	English	US			
Color	Rob Marsh	448	136	252	1000	Sam Clafir	40000	2.41E+08	Action Ad Johnny De Pirates of t	370704	54083	Stephen Gi	4	blackbear http://ww		484	English	US			
Color	Barry Sonr	451	106	188	718	Michael St	10000	1.79E+08	Action Ad Will Smith Men in Bla	268154	12572	Nicole Sch	1	alien crim	http://ww	341	English	US			
Color	Peter Jack	422	164	0	773	Adam Brov	5000	2.55E+08	Adventure Aidan Turn The Hobbi	354228	9152	James Nes	0	army elf j	http://ww	802	English	Ne			
Color	Marc Web	599	153	464	963	Andrew Ge	15000	2.62E+08	Action Ad Emma Sto The Amazi	451803	28489	Chris Zylka	0	lizard out	http://ww	1225	English	US			
Color	Ridley Sco'	343	156	0	738	William Hu	891	1.05E+08	Action Ad Mark Addy Robin Hoo	211765	3244	Scott Grin	0	1190s arc	http://ww	546	English	US			
Color	Peter Jack	509	186	0	773	Adam Brov	5000	2.58E+08	Adventure Aidan Turn The Hobbi	483540	9152	James Nes	6	dwarf elf j	http://ww	951	English	US			
Color	Chris Weit	251	113	129	1000	Eva Green	16000	70083519	Adventure Christoph The Golde	149019	24106	Kristin Sco	2	children e	http://ww	666	English	US			
Color	Peter Jack	446	201	0	84	Thomas Kr	6000	2.18E+08	Action Ad Naomi Wa King Kong	316018	7123	Evan Park	0	animal nar	http://ww	2618	English	Ne			
Color	James Can	315	194	0	794	Kate Winsl	29000	6.59E+08	Drama Ro Leonardo Titanic	793059	45223	Gloria Stu	0	artist love	http://ww	2528	English	US			
Color	Anthony R	516	147	94	11000	Scarlett Jo	21000	4.07E+08	Action Ad Robert Do Captain Ar	272670	64798	Chris Evan	0	based on c	http://ww	1022	English	US			
Color	Peter Berg	377	131	532	627	Alexander	14000	65173160	Action Ad Liam Nees Battleship	202382	26679	Tadanobu	0	box office	http://ww	751	English	US			
Color	Colin Trew	644	124	365	1000	Judy Greer	3000	6.52E+08	Action Ad Bryce Dallas Jurassic W	418214	8458	Omar Sy	0	dinosaur chil	http://ww	1290	English	US			
Color	Sam Mend	750	143	0	393	Helen McC	883	3.04E+08	Action Ad Albert Finni Skyfall	522030	2039	Rony Kinne	0	brawl chil	http://ww	1498	English	UK			
Color	Sam Raimi	300	135	0	4000	James Frar	24000	3.73E+08	Action Ad J.K. Simmc Spider-Ma	411164	43388	Kirsten Du	1	death doc	http://ww	1303	English	US			

45 DUPLICATE VALUES WERE FOUND AND REMOVED.
NEXT, I SELECTED ALL THE BLANK CELLS AND DELETED
ALL THE ROWS.

AutoSave IMDB_Movies humaira khan

File Home Insert Page Layout Formulas Data Review View Help

A1 color

color	director_n	num_critic	duration	director_f	actor_3	f	actor_2	n	actor_1	gross	genres	actor_1_n	movie	titl	num	vote	cast	total	actor_3_n	facenumb	plot	keyw	movie	imc	num	user	language	country	content				
2	Color	James Can	723	178	0	855	Joel David	1000	7.61E+08	Action Ad CCH Pounder	Avatar	886204	4834	Wes Studi	0	avatar fut	http://ww	3054	English	USA	PG-13												
3	Color	Gore Verbi	302	169	563	1000	Orlando Bl	40000	3.09E+08	Action Ad Johnny De Pirates of	Titanic	471220	48350	Jack Dauer	0	goddess n	http://ww	1238	English	USA	PG-13												
4	Color	Sam Mend	602	148	0	161	Rory Kinne	11000	2E+08	Action Ad Christoph	'Spectre'	275868	11700	Stephanie	1	bomb esp	http://ww	994	English	UK	PG-13												
5	Color	Christophe	813	164	22000	23000	Christian B	27000	4.48E+08	Action Ad Tom Hardy The Dark	K	1144337	106759	Joseph Go	0	deception	http://ww	2701	English	USA	PG-13												
6	Color	Doug Walker		131		Rob Walk	131			Document Ad Go To Special	?	X				0																	
7	Color	Andrew St.	462	132	475	530	Samantha	640	73058679	Ad	Walk Star Wars:	8		Polly Walk	1	alien ame	http://ww	738	English	USA	PG-13												
8	Color	Sam Raimi	392	156	0	4000	James Frar	24000	3.37E+08	Action Ad	Spider-Man	383056	46055	Kirsten Dui	0	sandman j	http://ww	1902	English	USA	PG-13												
9	Color	Nathan Gr	324	100	15	284	Donna Mu	799	2.01E+08	Action Ad	Avatar	294810	3244	M.C. Gain	1	17th centu	http://ww	387	English	USA	PG												
10	Color	Joss Whed	635	141	0	19000	Robert Do	26000	4.59E+08	Action Ad	Johnny De Pirates of	the Caribbean	471220	48350	Scarlett Jo	4	artificial	http://ww	1117	English	USA	PG-13											
11	Color	David Yate	375	153	282	10000	Daniel Rad	25000	3.02E+08	Action Ad	Christoph	'Spectre'	275868	11700	Rupert Gri	3	blood boc	http://ww	973	English	UK	PG											
12	Color	Zack Snyde	673	183	0	2000	Lauren Col	15000	3.3E+08	Action Ad	Tom Hardy The Dark	K	1144337	106759	Alan D. Pui	0	based on c	http://ww	3018	English	USA	PG-13											
13	Color	Bryan Sing	434	169	0	903	Marlon Br	18000	2E+08	Action Ad	Spider-Man	383056	46055	Frank Lang	0	crystal epi	http://ww	2367	English	USA	PG-13												
14	Color	Marc Forst	403	106	395	393	Mathieu A	451	1.68E+08	Action Ad	Avatar	294810	3244	Rory Kinne	1	action her	http://ww	1243	English	UK	PG-13												
15	Color	Gore Verbi	313	151	563	1000	Orlando Bl	40000	4.23E+08	Action Ad	Spider-Man	383056	46055	Jack Dauer	2	box office	http://ww	1832	English	USA	PG-13												
16	Color	Gore Verbi	450	150	563	1000	Ruth Wilsc	40000	89289910	Action Ad	Avatar	294810	3244	Tom Wilki	1	horse out	http://ww	711	English	USA	PG-13												
17	Color	Zack Snyde	733	143	0	748	Christophe	15000	2.91E+08	Action Ad	Tom Hardy The Dark	K	1144337	106759	Harry Len	0	based on c	http://ww	2536	English	USA	PG-13											
18	Color	Andrew Ac	258	150	80	201	Pierfrance	22000	1.42E+08	Action Ad	Spider-Man	383056	46055	Damijan A	4	brother br	http://ww	438	English	USA	PG												
19	Color	Joss Whed	703	173	0	19000	Robert Do	26000	4.59E+08	Action Ad	Johnny De Pirates of	the Caribbean	1144337	106759	Scarlett Jo	3	alien invas	http://ww	1722	English	USA	PG-13											
20	Color	Rob Mars	448	136	252	1000	Sam Claffir	40000	3.02E+08	Action Ad	Spider-Man	383056	46055	Peter Jack	0	blackbear c	http://ww	1290	English	USA	PG-13												
21	Color	Barry Sonn	451	106	188	718	Michael St	1000	4.07E+08	Action Ad	Robert Do Captain Ar	272670	64798	Chris Evan	0	dinosaur c	http://ww	2023	Rory Kinne	0	death doc	http://ww	1498	English	UK	PG-13							
22	Color	Peter Jack	422	164	0	773	Adam Bro	5000	2.58E+08	Action Ad	Aidan Turn The Hobbit	202382	26679	Tadanobu	0	box office	http://ww	751	English	USA	PG-13												
23	Color	Marc Web	599	153	464	963	Andrew Ge	15000	2.62E+08	Action Ad	Naomi Sto Amazi	202382	26679	Damijan A	2	children e	http://ww	1225	English	USA	PG-13												
24	Color	Ridley Sco	343	156	0	738	William Hu	891	1.05E+08	Action Ad	Mark Addy Robin Hoo	211765	3244	Tom Wilki	0	lizard out	http://ww	802	English	New Zeala	PG-13												
25	Color	Peter Jack	509	186	0	773	Adam Bro	5000	2.58E+08	Action Ad	Naomi Sto Amazi	202382	26679	Scarlett Jo	0	army eff	http://ww	951	English	USA	PG-13												
26	Color	Chris Weit	251	113	129	1000	Eva Green	16000	70083519	Action Ad	Spider-Man	383056	46055	Peter Jack	0	blackbear c	http://ww	2367	English	USA	PG-13												
27	Color	Peter Jack	446	201	0	84	Thomas Kr	6000	2.18E+08	Action Ad	Naomi Sto Amazi	202382	26679	Scarlett Jo	1	blackbear c	http://ww	1243	English	UK	PG-13												
28	Color	James Can	315	194	0	794	Kate Winsl	29000	6.59E+08	Action Ad	Avatar	294810	3244	Peter Jack	0	blackbear c	http://ww	2528	English	USA	PG-13												
29	Color	Anthony R	516	147	94	11000	Scarlett Jo	21000	4.07E+08	Action Ad	Robert Do Captain Ar	272670	64798	Chris Evan	0	based on c	http://ww	1022	English	USA	PG-13												
30	Color	Peter Berg	377	131	532	627	Alexander	14000	65173160	Action Ad	Aidan Turn The Hobbit	202382	26679	Scarlett Jo	0	box office	http://ww	751	English	USA	PG-13												
31	Color	Colin Trev	644	124	365	1000	Judy Greer	3000	6.52E+08	Action Ad	Brad Pitt Dali Jurassic W	202382	26679	Scarlett Jo	0	blackbear c	http://ww	1225	English	USA	PG-13												
32	Color	Sam Mend	750	143	0	393	Helen McC	883	3.04E+08	Action Ad	Albert Finni Skyfall	522030	2039	Rory Kinne	0	brawl chil	http://ww	666	English	UK	PG-13												
33	Color	Sam Raimi	300	135	0	4000	James Frar	24000	3.73E+08	Action Ad	J.K. Simmc Spider-Ma	411164	43388	Kirsten Dui	1	death doc	http://ww	2618	English	New Zeala	PG-13												

AutoSave IMDB_Movies humaira khan

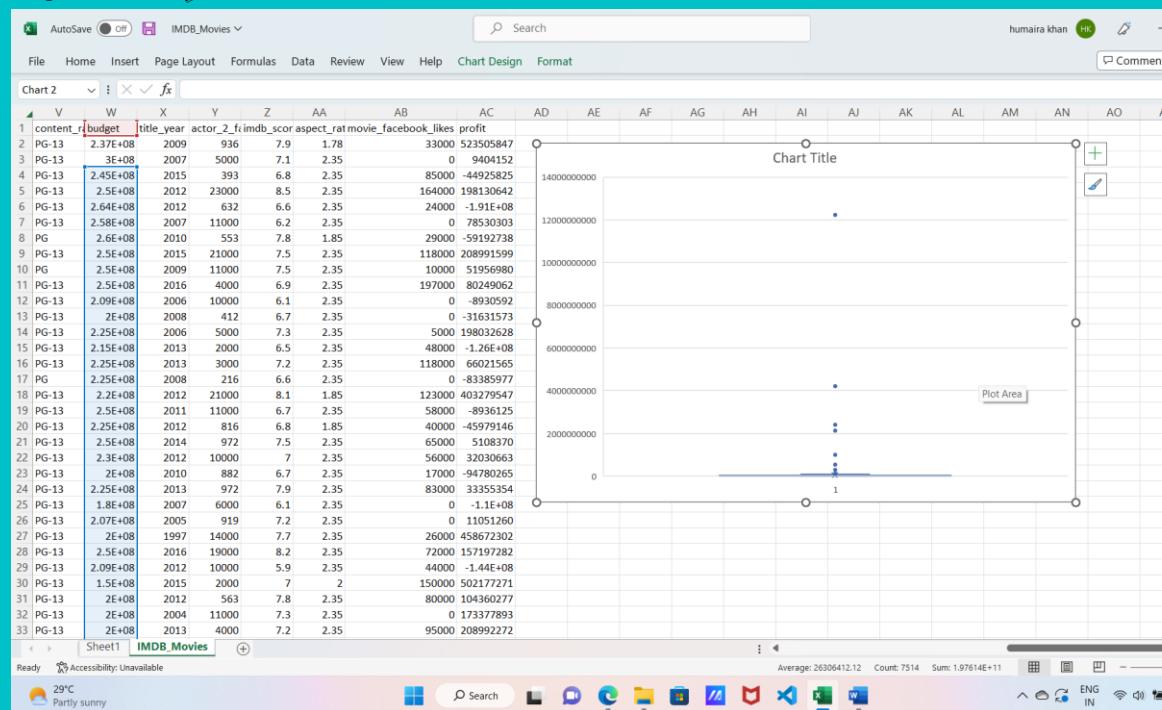
File Home Insert Page Layout Formulas Data Review View Help

A6 color

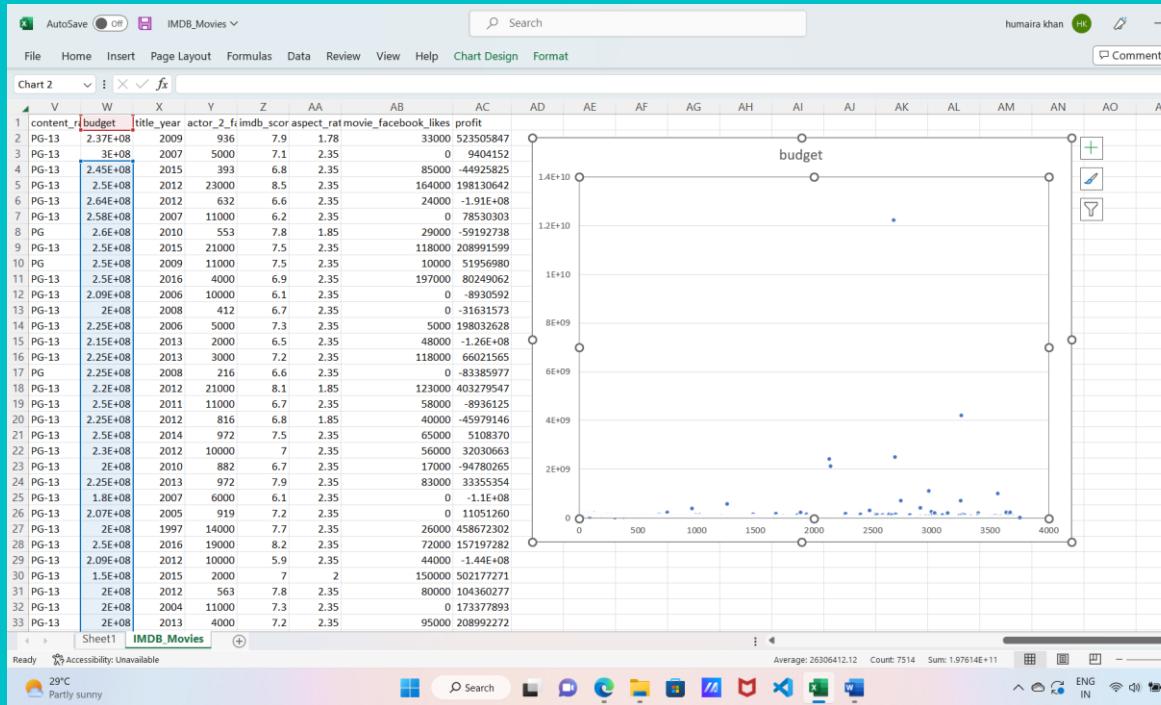
color	director_n	num_critic	duration	director_f	actor_3	f	actor_2	n	actor_1	gross	genres	actor_1_n	movie	titl	num	vote	cast	total	actor_3_n	facenumb	plot	keyw	movie	imc	num	user	language	country	content
2	Color	James Can	723	178	0	855	Joel David	1000	7.61E+08	Action Ad CCH Pounder	Avatar	886204	4834	Wes Studi	0	avatar fut	http://ww	3054	English	USA	PG-13								
3	Color	Gore Verbi	302	169	563	1000	Orlando Bl	40000	3.09E+08	Action Ad Johnny De Pirates of	Titanic	471220	48350	Jack Dauer	0	goddess n	http://ww	1238	English	USA	PG-13								
4	Color	Sam Mend	602	148	0	161	Rory Kinne	11000	2E+08	Action Ad Christoph	'Spectre'	275868	11700	Stephanie	1	bomb esp	http://ww	994	English	UK	PG-13								
5	Color	Christophe	813	164	22000	23000	Christian B	27000	4.48E+08	Action Ad Tom Hardy The Dark	K	1144337	106759	Joseph Go	0	deception	http://ww	2701	English	USA	PG-13								
6	Color	Doug Walker		131		Rob Walk	131			Document Ad Go To Special	?	X			Polly Walk	1	alien ame	http://ww	738	English	USA	PG-13							
7	Color	Andrew St.	462	132	475	530	Samantha	640	73058679	Ad	Walk Star Wars:	8		Kirsten Dui	0	sandman j	http://ww	1902	English	USA	PG-13								
8	Color	Sam Raimi	392	156	0	4000	James Frar	24000	3.37E+08	Action Ad	J.K. Simmc Spider-Ma	383056	46055	M.C. Gain	1	17th centu	http://ww	387	English	USA	PG								
9	Color	Nathan Gr	324	100	15	284	Donna Mu	799	2.01E+08	Action Ad	Brad Garre Tangled	294810	3244	Scarlett Jo	4	artificial	http://ww	1117	English	USA	PG-13								
10	Color	Joss Whed	635	141	0	19000	Robert Do	26000	4.59E+08	Action Ad	Chris Hem Avenger	202382	26679	Rupert Gri	2	blood boc	http://ww	973	English	UK	PG								
11	Color	David Yate	375	153	282	10000	Daniel Rad	25000	3.02E+08	Action Ad	Henry Cav Batman v	212904	1873	Polly Walk	1	alien ame	http://ww	1243	English	USA	PG-13								
12	Color	Zack Snyde	673	183	0	2000	Lauren Col	15000	3.3E+08	Action Ad	Henry Cav Batman v	212904	1873	James Nes	0	blackbear c	http://ww	3018	English	USA	PG-13								
13	Color	Bryan Sing	434	169	0	903	Marlon Br	18000	Delete	? X	Spac Superman	212904	1873	Scarlett Jo	0	blackbear c	http://ww	2367	English	USA	PG-13								
14	Color	Marc Forst	403	106	395	393	Mathieu A	451	40000	Delete	Quantum	303784	3244	Frank Lang	0	blackbear c	http://ww	1243	English	USA	PG-13								
15	Color	Gore Verbi	313	151	563	1000	Orlando Bl	40000	4.07E+08	Action Ad	De Peires of	25000	522040	Jack Dauer	0	blackbear c	http://ww	1832	English	USA	PG-13								
16	Color	Gore Verbi	450	150	563	1000	Ruth Wilsc	40000	0	Shift cells left	De Peires of	211765	3244	Harry Len	0	blackbear c	http://ww	711	English	USA	PG-13								
17	Color	Zack Snyde	733	143	0	748	Christophe	15000	0	Shift cells right	De Peires of	211765	3244	Damijan A	0	blackbear c													

I FIRST CREATED A NEW COLUMN PROFIT WHICH IS THE DIFFERENCE BETWEEN GROSS AND BUDGET. THEN I USED THE SORT FUNCTION ON PROFIT COLUMN(ASC TO DESCENDING)

NOW TO PLOT BUDGET AND PROFIT DATA, AND TO FIND OUTLIERS, THE BEST CHART WOULD BE BOXPLOT AND SCATTERPLOT.



WE CAN CLEARLY SEE THE OUTLIER HERE. BELOW IS SCATTER PLOT.



NOW TO FIND THE MAXIMUM PROFIT,

A screenshot of a Microsoft Excel spreadsheet titled "IMDB_Movies". The formula bar shows the formula `=MAX(AC:AC)`. The spreadsheet contains the following data:

	AD	AE	AF	AG	AH	AI
2		MAX	523505847			
3		MIN	-12213298588			
4						
5						
6						
7						
8						
9						
10						
11						

**WE USED THE MAX FUNCTION AND WE GOT TO KNOW THAT AVATAR IS THE MOVIE WITH HIGHEST PROFIT.
TO FIND THE TOP 10 WE USED THE SORT FUNCTION.**

TOP 10 HIGHEST MOVIES	PROFIT
Avatar	523505847
Jurassic World	502177271
Titanic	458672302
Stars War:Episode IV-A New Hope	449935665
E.T. the Extra-Terrestrial	424449459
The Avengers	403279547
The Lion King	377783777
Stars War:Episode I- The Phantom Menace	359544677
The Dark Knight	348316061
The Hunger Games	329999255

**THE THIRD QUERY WAS TO CALCULATE TOP 250 MOVIES WITH THE HIGHEST IDMB RATING.
FOR THIS WE USED THE PIVOT TABLE FUNCTION, SORTED THE VALUES IN DESCENDING ORDER AND IT SHOWS THE TOP RATED MOVIES.**

The screenshot shows a Microsoft Excel spreadsheet titled "IMDB_Movies". The PivotTable Fields pane on the right indicates that the "imdb_score" field is selected. The main area displays a list of movies from highest to lowest IMDB score. The data includes:

Score	Movie Title
9.3	The Shawshank Redemption
9.2	The Godfather
9	The Dark Knight
9	The Godfather: Part II
8.9	Pulp Fiction
8.9	Schindler's List
8.9	The Good, the Bad and the Ugly
8.9	The Lord of the Rings: The Return of the King
8.8	Fight Club
8.8	Forrest Gump
8.8	Inception
8.7	Star Wars: Episode V - The Empire Strikes Back
8.7	The Lord of the Rings: The Fellowship of the Ring
8.7	City of God
8.6	Goodfella
8.6	One Flew Over the Cuckoo's Nest
8.6	Seven Samurai
8.6	Star Wars: Episode IV - A New Hope
8.6	The Lord of the Rings: The Two Towers
8.6	The Matrix
8.6	American History X
8.6	Interstellar
8.6	Modern Times

	A	B	C
34.	Spirited Away		
35.	The Silence of the Lambs		
36.	The Usual Suspects		
37.	Alien		
38.	Apocalypse Now		
39.	Back to the Future		
40.	Chungking Express		
41.	Django Unchained		
42.	Gladiator		
43.	Memento		
44.	Psycho		
45.	Raiders of the Lost Ark		
46.	Samsara		
47.	Titanic		
48.	Terminator 2: Judgment Day		
49.	The Dark Knight Rises		
50.	The Departed		
51.	The Green Mile		
52.	The Lion King		
53.	The Lives of Others		
54.	The Pianist		
55.	The Prestige		
56.	Whiplash		
57.			8.4
58.	A Separation		
59.	Aliens		
60.	Amélie		
61.	American Beauty		
62.	Braveheart		
63.	Das Boot		
64.	Lawrence of Arabia		
65.	Oldboy		
66.	Once Upon a Time in America		

	A	B	C	D
67.	Princess Mononoke			
68.	Requiem for a Dream			
69.	Reservoir Dogs			
70.	Star Wars: Episode VI - Return of the Jedi			
71.	WALL-E			
72.			8.3	
73.	2001: A Space Odyssey			
74.	Armageddon			
75.	Batman Begins			
76.	Downfall			
77.	Eternal Sunshine of the Spotless Mind			
78.	Good Will Hunting			
79.	Hoop Dreams			
80.	Indiana Jones and the Last Crusade			
81.	Inception			
82.	Inside Out			
83.	L.A. Confidential			
84.	Metropolis			
85.	Monty Python and the Holy Grail			
86.	Raging Bull			
87.	Room			
88.	Seven			
89.	Snatch			
90.	Some Like It Hot			
91.	The Hunt			
92.	The Sting			
93.	Toy Story			
94.	Toy Story 3			
95.	Unforgiven			
96.	Up			
97.			8.2	
98.	A Beautiful Mind			
99.	Blade Runner			

	A	B	C	D	E
97.		8.2			
98.	A Beautiful Mind				
99.	Blade Runner				
100.	Captain America: Civil War				
101.	Casino				
102.	Die Hard				
103.	Finding Nemo				
104.	Gone with the Wind				
105.	Gran Torino				
106.	How to Train Your Dragon				
107.	Howl's Moving Castle				
108.	Incendies				
109.	Into the Wild				
110.	Lock, Stock and Two Smoking Barrels				
111.	On the Waterfront				
112.	Pan's Labyrinth				
113.	The Act of Killing				
114.	The Big Lebowski				
115.	The Bridge on the River Kwai				
116.	The Secret in Their Eyes				
117.	The Thing				
118.	The Wolf of Wall Street				
119.	Trainspotting				
120.	V for Vendetta				
121.	Warrior				
122.		8.1			
123.	12 Years a Slave				
124.	Akira				
125.	Amores Perros				
126.	Annie Hall				
127.	Before Sunrise				
128.	Butch Cassidy and the Sundance Kid				
129.	Deadpool				

127	Before Sunrise	Ã			
128	Butch Cassidy and the Sundance Kid	Ã			
129	Deadpool	Ã			
130	Donnie Darko	Ã			
131	Elite Squad	Ã			
132	Gone Girl	Ã			
133	Groundhog Day	Ã			
134	Guardians of the Galaxy	Ã			
135	Hotel Rwanda	Ã			
136	In the Shadow of the Moon	Ã			
137	Jurassic Park	Ã			
138	Kill Bill: Vol. 1	Ã			
139	Mad Max: Fury Road	Ã			
140	Million Dollar Baby	Ã			
141	Monsters, Inc.	Ã			
142	No Country for Old Men	Ã			
143	Pirates of the Caribbean: The Curse of the Black Pearl	Ã			
144	Platoon	Ã			
145	Prisoners	Ã			
146	Rocky	Ã			
147	Rush	Ã			
148	Shutter Island	Ã			
149	Sin City	Ã			
150	Spotlight	Ã			
151	Stand by Me	Ã			
152	Tae Guk Gi: The Brotherhood of War	Ã			
153	The Avengers	Ã			
154	The Best Years of Our Lives	Ã			
155	The Bourne Ultimatum	Ã			
156	The Celebration	Ã			
157	The Grand Budapest Hotel	Ã			
158	The Help	Ã			
159	The Imitation Game	Ã			
		IMDB_Movies			
		Ready	Accessibility: Unavailable		
		27°C	Haze		

160	The Martian	Ã			
161	The Princess Bride	Ã			
162	The Revenant	Ã			
163	The Sea Inside	Ã			
164	The Sixth Sense	Ã			
165	The Terminator	Ã			
166	The Truman Show	Ã			
167	The Wizard of Oz	Ã			
168	There Will Be Blood	Ã			
169	Woodstock	Ã			
170	8				
171	A Fistful of Dollars	Ã			
172	Aladdin	Ã			
173	Before Sunset	Ã			
174	Big Fish	Ã			
175	Black Swan	Ã			
176	Blood Diamond	Ã			
177	Blood In, Blood Out	Ã			
178	Bowling for Columbine	Ã			
179	Boyhood	Ã			
180	Brazil	Ã			
181	Casino Royale	Ã			
182	Catch Me If You Can	Ã			
183	Central Station	Ã			
184	Cinderella Man	Ã			
185	Dallas Buyers Club	Ã			
186	Dancer in the Dark	Ã			
187	Dances with Wolves	Ã			
188	Dead Poets Society	Ã			
189	District 9	Ã			
190	Doctor Zhivago	Ã			
191	Fiddler on the Roof	Ã			
192	Her	Ã			
		IMDB_Movies			
		Ready	Accessibility: Unavailable		
		27°C	Haze		

193	In Bruges	A	B	C	D	A250	
194	Jaws						The Sound of Music
195	JFK						The Straight Story
196	Kill Bill: Vol. 2						True Romance
197	Life of Pi						Waltz with Bashir
198	Magnolia						X-Men: Days of Future Past
199	Mulholland Drive						Young Frankenstein
200	My Name Is Khan						7.9
201	Mystic River						4 Months, 3 Weeks and 2 Days
202	Persepolis						Almost Famous
203	Rain Man						Amour
204	Ratatouille						Avatar
205	Serenity						Before Midnight
206	Shaun of the Dead						Big Hero 6
207	Sicko						Boogie Nights
208	Sling Blade						Captain Phillips
209	Slumdog Millionaire						Children of Men
210	Star Trek						Crash
211	The Artist						Crouching Tiger, Hidden Dragon
212	The Exorcist						Do the Right Thing
213	The Incredibles						E.T. the Extra-Terrestrial
214	The Iron Giant						Ed Wood
215	The King's Speech						Edge of Tomorrow
216	The Perks of Being a Wallflower						Edward Scissorhands
217	The Pursuit of Happyness						Ernest & Celestine
218	The Sound of Music						Glory
219	The Straight Story						Halloween
220	True Romance						Hero
221	Waltz with Bashir						Hot Fuzz
222	X-Men: Days of Future Past						How to Train Your Dragon 2
223	Young Frankenstein						Iron Man
224	7.9						Letters from Iwo Jima
225	4 Months, 3 Weeks and 2 Days						Little Miss Sunshine
							Moon

4) FORTH QUERY

NOW THE NEXT ONE WAS TO CALCULATE TOP 10 DIRECTORS WITH A MEAN OF HIGHEST RATING MOVIES. AGAIN WE USED THE PIVOT TABLE.

The screenshot shows a Microsoft Excel spreadsheet titled "IMDB_Movies1". The Pivot Table is set up with "Row Labels" as "Average of imdb_score". The data shows the average IMDB score for various directors. The "Grand Total" row at the bottom also displays the value 8.47. The table includes columns for the director's name and their average score.

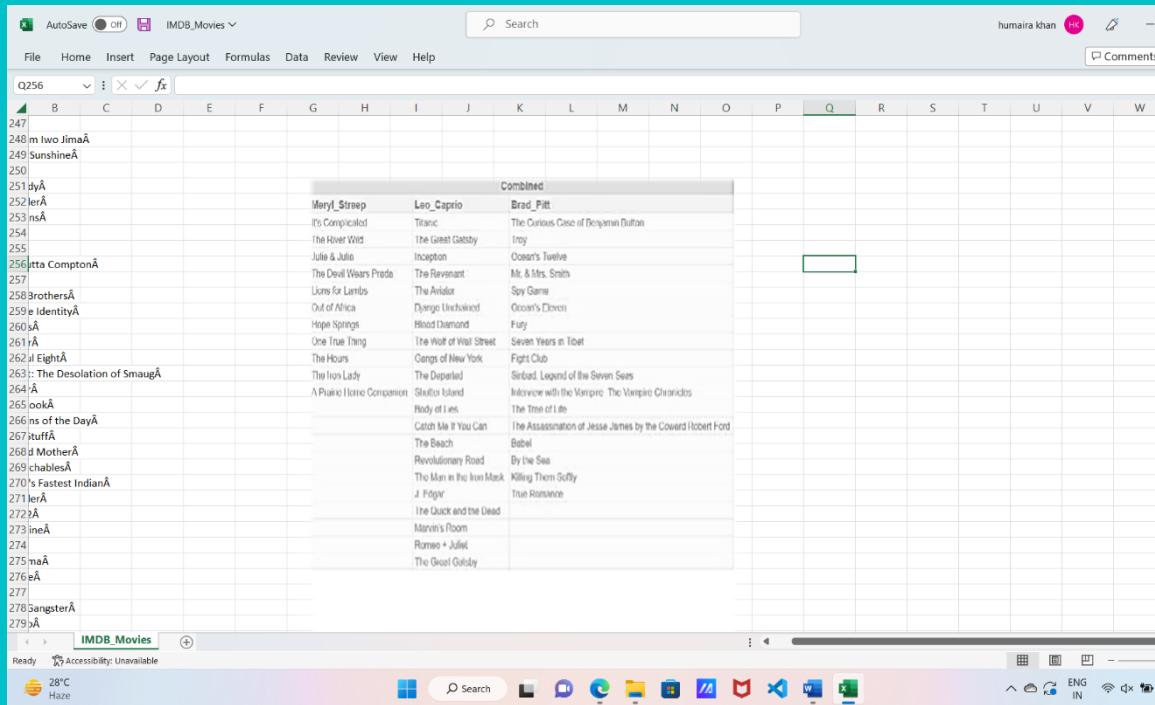
	A	B
2		
3	Row Labels	Average of imdb_score
4	Akira Kurosawa	8.7
5	Charles Chaplin	8.6
6	Tony Kaye	8.6
7	Alfred Hitchcock	8.5
8	Damien Chazelle	8.5
9	Majid Majidi	8.5
10	Ron Fricke	8.5
11	Sergio Leone	8.433333333
12	Christopher Nolan	8.425
13	Asghar Farhadi	8.4
14	Richard Marquand	8.4
15	Grand Total	8.47
16		

5) THE FIFTH QUERY WAS TO FIND THE POPULAR GENRES. WE DID THIS BY SEEING THE IDMB SCOREES. TOP 10 GENRES ARE:

The screenshot shows a Microsoft Excel spreadsheet titled "IMDB_Movies". The table has three columns: "Rank", "Top 10 popular genres", and "mean_imdb_score". The data is as follows:

Rank	Top 10 popular genres	mean_imdb_score
1	Adventure Animation Drama Family Musical	8.5
2	Crime Drama Fantasy Mystery	8.5
3	Action Adventure Drama Fantasy War	8.4
4	Adventure Animation Fantasy	8.4
5	Adventure Drama Thriller War	8.4
6	Adventure Animation Comedy Drama Family Fantasy	8.3
7	Biography Drama History Music	8.3
8	Documentary Drama Sport	8.3
9	Documentary War	8.3
10	Adventure Drama War	8.3

6) THE LAST QUERY WAS REALLY CHALLENGING.
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.



Meryl_Streep	Leo_Caprio	Brad_Pitt
It's Complicated	Titanic	The Curious Case of Benjamin Button
The Devil Wears Prada	The Great Gatsby	Troy
The Hours	Inception	Ocean's Twelve
The Iron Lady	The Revenant	Mr. & Mrs. Smith
The King's Speech	The Aviator	Spy Game
The Merchant of Venice	Django Unchained	Occident's Eleven
The Pianist	Hugo	Fury
The Reader	The Wolf of Wall Street	Seven Years in Tibet
The Silence of the Lambs	Gangs of New York	Fight Club
The Social Network	The Departed	Sir Bodin: Legend of the Seven Seas
The Thin Red Line	Shutter Island	Interview with the Vampire: The Vampire Chronicles
The Tree of Life	Moneyball	The Tree of Life
Catch Me If You Can	The Assassination of Jesse James by the Coward Robert Ford	Babel
The Beach	Revolutionary Road	By the Sea
The Man in the Iron Mask	The Man in the Iron Mask	Killing Them Softly
J. Edgar	J. Edgar	True Romance
The Quick and the Dead	Marvin's Room	
Marvin's Room	Romeo + Juliet	
Romeo + Juliet	The Great Gatsby	

FIND THE MEAN OF THE NUM_CRITIC_FOR_REVIEWS AND NUM_USERS_FOR REVIEW AND IDENTIFY THE ACTORS WHICH HAVE THE HIGHEST MEAN.

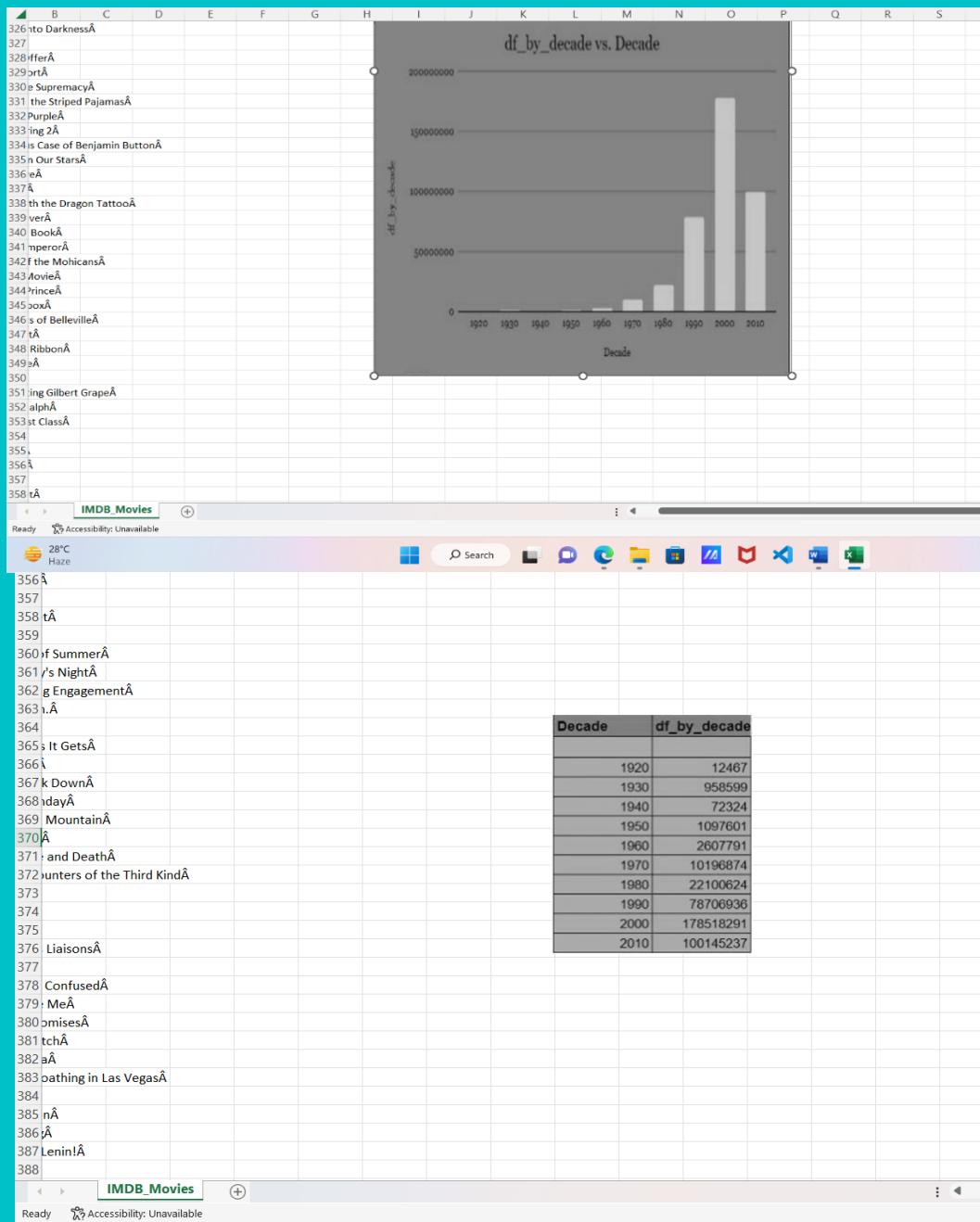
The screenshot shows a Microsoft Excel spreadsheet titled "IMDB_Movies". A table titled "Top 5 Actors who have the highest mean in num_critic_for_reviews" is displayed. The table has three columns: Rank, Actor, and mean num_critic_for_reviews. The data is as follows:

Rank	Actor	mean num_critic_for_reviews
1	Albert Finney	750.00
2	Phaldut Sharma	739.00
3	Peter Capaldi	654.00
4	Craig Stark	596.00
5	Bérénice Bejo	576.00

The screenshot shows a Microsoft Excel spreadsheet titled "IMDB_Movies". A table titled "Top 5 Actors who have the highest mean in num_user_for_reviews" is displayed. The table has three columns: Rank, Actor, and Mean. The data is as follows:

Rank	Actor	Mean
1	Heather Donahue	3400.00
2	Christo Jivkov	2814.00
3	Steve Bastoni	2789.00
4	Phaldut Sharma	1885.00
5	Orlando Bloom	1842.00

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.



EXPLORATORY DATA ANALYSIS FOR BANK LOAN

PROJECT DESCRIPTION - THE LOAN PROVIDING COMPANIES FIND IT HARD TO GIVE LOANS TO THE PEOPLE DUE TO THEIR INSUFFICIENT OR NON-EXISTENT CREDIT HISTORY. BECAUSE OF THAT, SOME CONSUMERS USE IT AS THEIR ADVANTAGE BY BECOMING A DEFULTER. WE HAVE TO USE EDA TO ANALYZE THE PATTERNS PRESENT IN THE DATA. THIS WILL ENSURE THAT THE APPLICANTS CAPABLE OF REPAYING THE LOAN ARE NOT REJECTED.

APPROACH - THIS CASE STUDY AIMS TO IDENTIFY PATTERNS WHICH INDICATE IF AN APPLICANT HAS DIFFICULTY IN PAYING HIS/HER 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.

I HAVE MAINLY FOCUSED ON ANALYZING PREVIOUS_APPLICATION.CSV I.E. DATA ABOUT PREVIOUS APPLICATION OF AN APPLICANT. FOR THE EXPLORATORY DATA ANALYSIS, MENTIONED

STEPS HAVE BEEN FOLLOWED. → IMPORT MODULES,

READ THE DATASET

SK_ID_PREV	SK_ID_CURR	NAME_CONTRACT_TYPE	AMT_ANNUITY	AMT_APPLICATION	AMT_CREDIT	AMT_DOWN_PAYMENT
2030495	271877	Consumer loans	1730.430	17145.0	17145.0	0.0
2802425	108129	Cash loans	25188.615	607500.0	679671.0	Nan
2523466	122040	Cash loans	15060.735	112500.0	136444.5	Nan
2819243	176158	Cash loans	47041.335	450000.0	470790.0	Nan
1784265	202054	Cash loans	31924.395	337500.0	404055.0	Nan

Data columns (total 37 columns):			
#	Column	Non-Null Count	Dtype
0	SK_ID_PREV	1670214	non-null
1	SK_ID_CURR	1670214	non-null
2	NAME_CONTRACT_TYPE	1670214	non-null
3	AMT_ANNUITY	1297979	non-null
4	AMT_APPLICATION	1670214	non-null
5	AMT_CREDIT	1670213	non-null
6	AMT_DOWN_PAYMENT	774370	non-null
7	AMT_GOODS_PRICE	1284699	non-null
8	WEEKDAY_APPR_PROCESS_START	1670214	non-null
9	HOUR_APPR_PROCESS_START	1670214	non-null
10	FLAG_LAST_APPL_PER_CONTRACT	1670214	non-null
11	NFLAG_LAST_APPL_IN_DAY	1670214	non-null
12	RATE_DOWN_PAYMENT	774370	non-null
13	RATE_INTEREST_PRIMARY	5951	non-null
14	RATE_INTEREST_PRIVILEGED	5951	non-null
15	NAME_CASH_LOAN_PURPOSE	1670214	non-null
16	NAME_CONTRACT_STATUS	1670214	non-null
17	DAYS_DECISION	1670214	non-null
18	NAME_PAYMENT_TYPE	1670214	non-null
19	CODE_REJECT_REASON	1670214	non-null
20	NAME_TYPE_SUITE	849889	non-null
21	NAME_CLIENT_TYPE	1670214	non-null
22	NAME_GOODS_CATEGORY	1670214	non-null
23	NAME_PORTFOLIO	1670214	non-null
24	NAME_PRODUCT_TYPE	1670214	non-null
25	CHANNEL_TYPE	1670214	non-null
26	SELLERPLACE_AREA	1670214	non-null
27	NAME_SELLER_INDUSTRY	1670214	non-null
28	CNT_PAYMENT	1297984	non-null
29	NAME_YIELD_GROUP	1670214	non-null
30	PRODUCT_COMBINATION	1669868	non-null
31	DAYS_FIRST_DRAWING	997149	non-null
32	DAYS_FIRST_DUE	997149	non-null
33	DAYS_LAST_DUE_1ST_VERSION	997149	non-null
34	DAYS_LAST_DUE	997149	non-null
35	DAYS_TERMINATION	997149	non-null
36	NFLAG_INSURED_ON_APPROVAL	997149	non-null

PREV_AP_DF CONTAINS 37 FEATURES AND 1670214 ROWS(OUT OF WHICH 15 FEATURES ARE FLOAT64, 6 FEATURES ARE INTEGER, 16 FEATURES ARE OBJECT DATATYPE)

FOLLOWING ARE THE COMMON FEATURES AMONG APPLICATION DATA AND PREVIOUS APPLICATION DATA -

['SK_ID_CURR', 'NAME_CONTRACT_TYPE', 'AMT_CREDIT', 'AMT_ANNUITY', 'AMT_GOODS_PRICE', 'NAME_TYPE_SUITE',

'WEEKDAY_APPR_PROCESS_START', 'HOUR_APPR_PROCESS_START']

SK_ID_CURR IS AN UNIQUE IDENTIFIER, WHICH WILL USE TO MERGE THE RELEVANT COLUMNS OF 2 DATAFRAMES (APPLICATION DATA AND PREVIOUS APPLICATION DATA).

DATA CLEANING

MISSING VALUE HANDLING, TYPE CASTING, FIXING ROWS AND COLUMNS – REMOVING UNNECESSARY ROWS/COLUMNS (THROUGH MISSING VALUE HANDLING AND CORRELATION), HANDLING OUTLIERS.

FIRST I HAVE CALCULATED THE MISSING VALUE PERCENTAGES FOR EACH FEATURE IN PREVIOUS APPLICATION DATA.

	category	percentage
5	RATE_INTEREST_PRIMARY	99.843698
6	RATE_INTEREST_PRIVILEGED	99.843698
2	AMT_DOWN_PAYMENT	53.636480
4	RATE_DOWN_PAYMENT	53.636480
7	NAME_TYPE_SUITE	49.119754
10	DAYS_FIRST_DRAWING	40.298129
11	DAYS_FIRST_DUE	40.298129
12	DAYS_LAST_DUE_1ST_VERSION	40.298129
13	DAYS_LAST_DUE	40.298129
14	DAYS_TERMINATION	40.298129
15	NFLAG_INSURED_ON_APPROVAL	40.298129
3	AMT_GOODS_PRICE	23.081773
0	AMT_ANNUITY	22.288665
8	CNT_PAYMENT	22.286366
9	PRODUCT_COMBINATION	0.020718
1	AMT_CREDIT	0.000060

THERE ARE 16 FEATURES IN PREV_APP_DF THAT HAVE MISSING VALUES.

PERMANENTLY DROPPING THE FEATURES

(RATE_INTEREST_PRIMARY AND RATE_INTEREST_PRIVILEGED) AS 99% DATA IS MISSING.

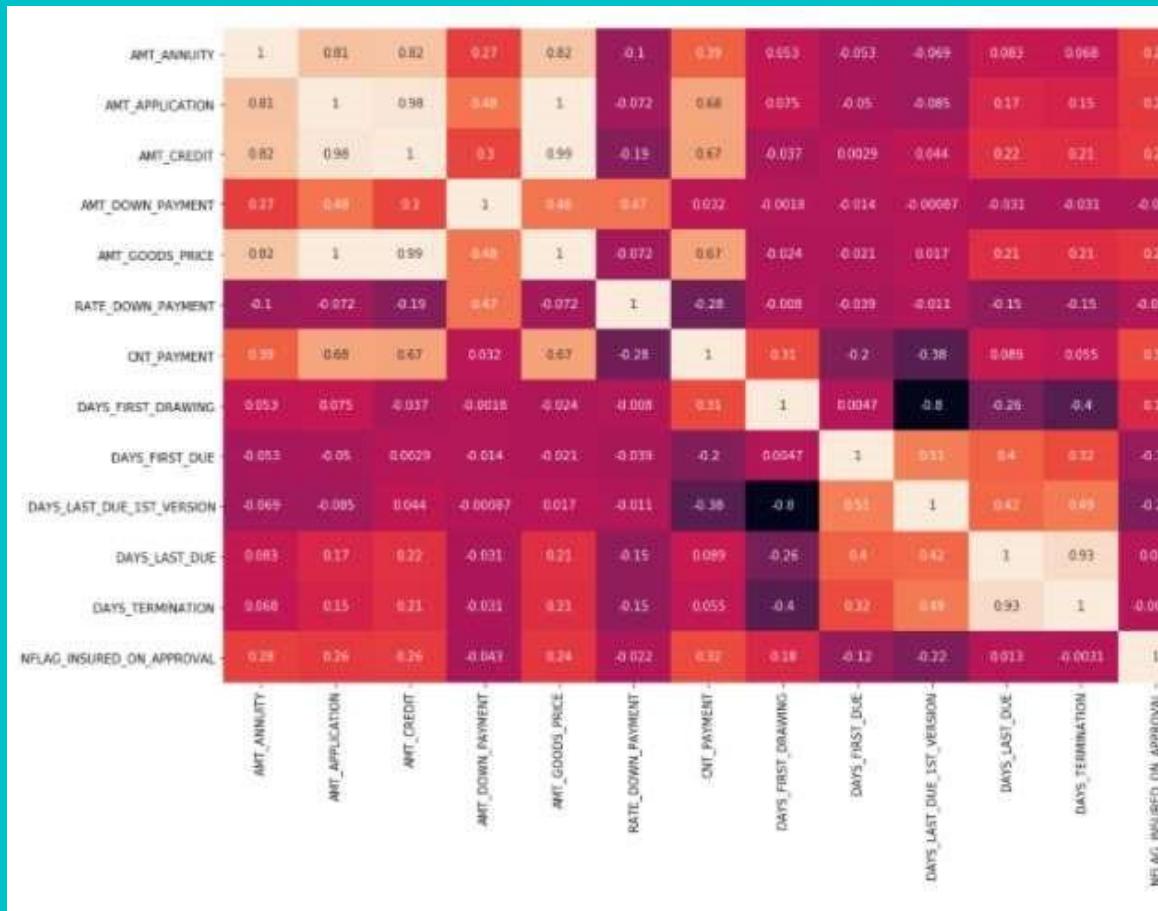
DROPPING ROWS CONTAINING MISSING VALUES FOR THE FEATURES (AMT_CREDIT AND PRODUCT_COMBINATION) FOR VERY LOW % OF MISSING DATA. DROPPING

ENTRIES WOULD NOT CAUSE IMPACT THE ANALYSIS AS PERCENTAGE OF MISSING VALUE IS VERY LOW (~2%).

→

UNIVARIATE ANALYSIS, BIVARIATE AND MULTIVARIATE ANALYSIS

FIRST, I EXTRACTED THE NUMERICAL VARIABLES IN THE DATASET AND CHECKED OUT THE CORRELATION COEFFICIENTS WITH THE HELP OF A HEATMAP.



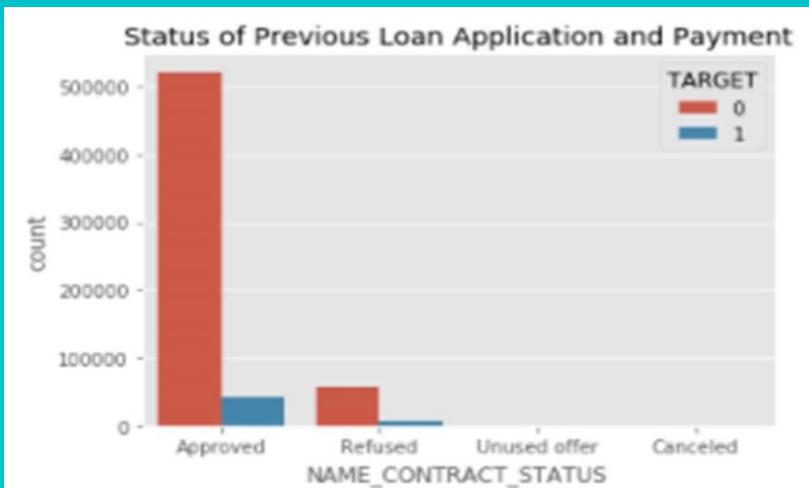
'DAYS_LAST_DUE' AND 'DAYS_TERMINATION' ARE HIGHLY CORRELATED
 'DAYS_FIRST_DRAWING' AND 'DAYS_LAST_DUE_1ST_VERSION' HAVE HIGH NEGATIVE CORRELATION

'AMT_ANNUITY','AMT_APPLICATION','AMT_CREDIT','AMT_GOODS_PRICE' ARE HIGHLY CORRELATED

THE FEATURES CAN BE REMOVED BEFORE MODELLING THIS DATA, AS THEY WOULD CAUSE COLLINEARITY

'DAYS_TERMINATION','DAYS_LAST_DUE_1ST_VERSION','AMT_APPLICATION','AMT_CREDIT','AMT_GOODS_PRICE' FOR EDA PURPOSE WE ARE NOT REMOVING THEM.

CHECKING DATA IMBALANCE IN PREVIOUS APPLICATION DATA



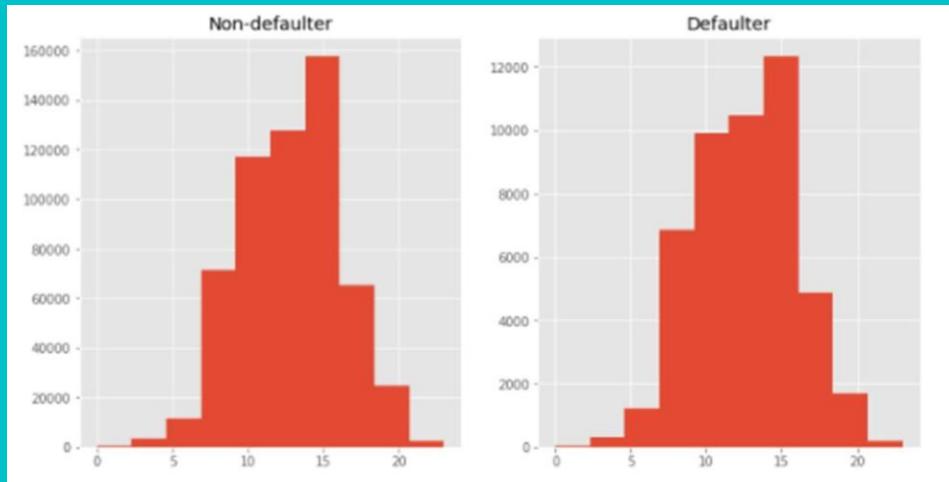
THE APPLICANTS WHOSE PREVIOUS LOANS WERE APPROVED ARE MORE LIKELY TO PAY CURRENT LOAN IN TIME, THAN THE APPLICANTS WHOSE PREVIOUS LOANS WERE REJECTED.

7% OF THE PREVIOUSLY APPROVED LOAN APPLICANTS THAT DEFULTED IN CURRENT LOAN

90 % OF THE PREVIOUSLY REFUSED LOAN APPLICANTS THAT WERE ABLE TO PAY CURRENT LOAN

THIS DATA IS HIGHLY IMBALANCED AS NUMBER OF DEFULTERS IS VERY LESS IN TOTAL POPULATION.

ANALYSIS OF NUMERIC FEATURES OF PREVIOUS APPLICATION DATA



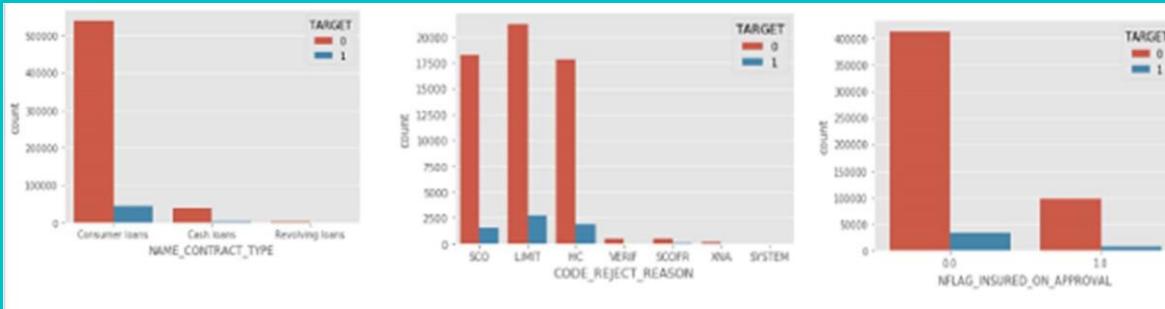
NUMBER
OF

DEFULTERS ARE LESS FOR LARGER AMOUNT OF ANNUITY OF PREVIOUS APPLICATION.

FOR HIGHER DOWN PAYMENT, DEFULTER CASES ARE LESS.

MOST OF THE LOANS ARE APPLIED AROUND 15:00 HOURS. THIS FEATURE IS DOES NOT HAVE VISIBLE IMPACT ON TARGET VARIABLE.

**ANALYSIS OF CATEGORICAL FEATURES OF
PREVIOUS APPLICATION DATA HIGHEST
NUMBER OF LOANS ARE APPLIED FOR
CONSUMER LOANS.**



AS SEEN IN THE ABOVE PLOT, 'SCO', 'LIMIT' AND 'HC' ARE THE MOST COMMON REASON OF REJECTION.

MOST OF THE PEOPLE DID NOT REQUEST INSURANCE DURING PREVIOUS LOAN APPLICATION.

MOST OF THE APPLICANTS ARE REPEATER.

'CASH THROUGH THE BANK' IS THE MOST

FREQUENTLY USED PAYMENT METHOD. THEN FOR A

GIVEN CATEGORICAL FEATURE, I OBTAINED

PERCENTAGE OF DEFAULTERS.

		Value	Percentage of Defaulter
23	Insurance	10.526316	
0	Vehicles	10.267410	
14	Jewelry	9.124951	
17	Auto Accessories	9.029763	
3	Mobile	8.615336	
15	Office Appliances	8.307692	
8	Computers	8.074335	
20	Weapon	8.064516	
21	Direct Sales	8.024691	
5	Audio/Video	7.698706	
7	Photo / Cinema Equipment	7.455000	
18	Sport and Leisure	7.354150	
2	Consumer Electronics	7.068548	
4	Construction Materials	6.978320	
9	XNA	6.885879	
24	Additional Service	6.730769	
6	Gardening	6.723063	
11	Homewares	6.708444	
19	Medicine	6.198747	
25	Education	5.882353	
1	Furniture	5.860781	
10	Clothing and Accessories	5.807427	
13	Other	5.765921	
12	Medical Supplies	5.584190	
16	Tourism	4.444444	
22	Fitness	4.268293	
26	Animals	0.000000	

Highest percentage of default cases are for i

	Value	Percentage of Defaulter
1	walk-in	9.165550
0	XNA	7.686995
2	x-sell	6.036420

highest.

	Value	Percentage of Defaulter
1	walk-in	9.165550
0	XNA	7.686995
2	x-sell	6.036420

From all the walk-in applicants 9% defaulted on loan.

	Value	Percentage of Defaulter
4	AP+ (Cash loan)	15.000000
1	Country-wide	7.908171
2	Regional / Local	7.551291
0	Stone	7.294692
3	Credit and cash offices	6.124197
5	Contact center	4.545455
6	Car dealer	0.000000

WHO PREVIOUSLY APPLIED FOR INSURANCE AND VEHICLES.

FOR CARDS DEFULTER RATE 15% LOAN APPLICANTS DEFULTED FOR

		Value	Percentage of Defaulter
0	Auto technology	10.522088	
9	Jewelry	9.019221	
3	Connectivity	8.780637	
2	Consumer electronics	7.451983	
7	Industry	7.211664	
4	Construction	6.597424	
5	XNA	6.226598	
1	Furniture	5.924492	
6	Clothing	5.857399	
8	Tourism	4.778157	
10	MLM partners	4.654655	

AP + CASH LOAN

IN SELLER INDUSTRY “AUTO TECHNOLOGY” HAS HIGHEST RATE OF DEFULTER

MLM PARTNERS HAS LOWEST NUMBER OF DEFULTERS

		Value	Percentage of Defaulter
4	XNA	17.119695	
2	high	8.340935	
1	middle	7.558098	
0	low_normal	6.844973	
3	low_action	6.608936	

DEFALTER PERCENTAGE IS HIGHEST WHERE NAME_YIELD_GROUP IS NOT KNOWN.
HIGHEST PERCENTAGE OF DEFAULT CASES IS FOR CARD STREET.

		Value	Percentage of Defaulter
13	Card Street	17.195005	
4	POS mobile with interest	8.781056	
0	POS other with interest	7.953141	
3	POS mobile without interest	7.888514	
2	POS household with interest	7.752151	
11	POS others without interest	7.256127	
15	Card X-Sell	6.888867	
5	POS household without interest	6.649376	
9	Cash Street: middle	6.475391	
10	Cash Street: high	6.417625	
8	Cash X-Sell: high	6.410114	
1	POS industry with interest	6.350635	
12	Cash X-Sell: middle	6.017039	
7	Cash Street: low	5.978876	
6	POS industry without interest	4.711940	
14	Cash X-Sell: low	3.988711	

SUMMARY OF PREVIOUS APPLICATION DATA

THERE ARE FEATURE COLUMNS IN THE DATASET THAT ARE HIGHLY CORRELATED TO EACH OTHER. WHICH MEANS BOTH WILL HAVE SIMILAR IMPACT ON THE TARGET VALUE. THOSE FEATURES CAN BE REMOVED BEFORE FEEDING THIS DATA TO A MODEL TO AVOID COLLINEARITY. FEATURE COLUMNS WITH 50% OR MORE MISSING DATA CAN BE DROPPED. THIS DATASET IS HIGHLY IMBALANCED

THE APPLICANTS WHOSE PREVIOUS LOANS WERE APPROVED ARE MORE LIKELY TO PAY CURRENT LOAN IN TIME, THAN THE APPLICANTS WHOSE PREVIOUS LOANS WERE REJECTED. NAME_CONTRACT_STATUS IS AN IMPORTANT FEATURE.

7% OF THE PREVIOUSLY APPROVED LOAN APPLICANTS

THAT DEFAULTED IN CURRENT LOAN 90 % OF THE

PREVIOUSLY REFUSED LOAN APPLICANTS THAT WERE

ABLE TO PAY CURRENT LOAN.

'SCO', 'LIMIT' AND 'HC' ARE THE MOST COMMON REASON OF REJECTION.

MOST OF THE PEOPLE DID NOT REQUEST INSURANCE DURING PREVIOUS LOAN APPLICATION.

FOR "CARDS" DEFaulTER PERCENTAGE IS HIGHEST (17%). 15% LOAN APPLICANT DEFAULTED FOR AP+ (CASH LOAN).

CURRENT APPLICATION DATA ANALYSIS:

FOR ANALYZING CURRENT APPLICATION DATA, I HAVE TAKEN A DIFFERENT APPROACH THAN THAT OF PREVIOUS APPLICATION DATA. BY TAKING A CLOSE LOOK AT THE FEATURES, I COULD IDENTIFY THE FEATURES OF DIFFERENT ASPECTS OF THE LOAN APPLICANT.

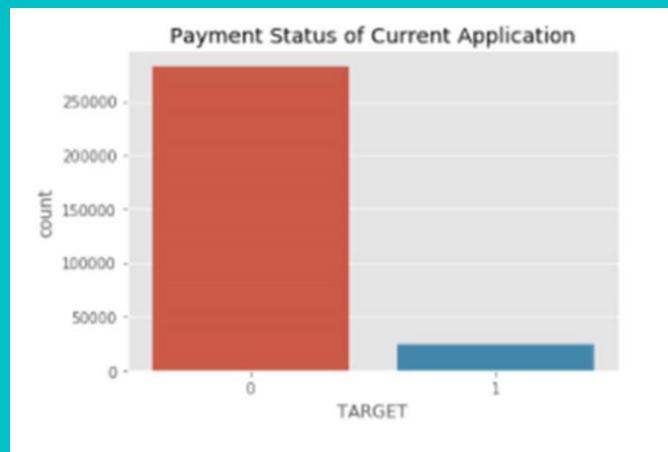
THAT'S WHY I HAVE DIVIDED THE FEATURES INTO SMALL SEGMENTS AND ANALYZED THEM SEGMENT-WISE USING A SMALLER DATA-FRAME CONTAINING ONLY RELEVANT CATEGORIES.

DATA CLEANING, MISSING DATA HANDLING, TYPE CASTING ARE DONE SEGMENT-WISE.

PLOTS AND PERCENTAGE WISE DEFULTER CALCULATION ARE DONE SEGMENT-WISE AS WELL.

CHECKING DATA IMBALANCE:

THIS DATA IS HIGHLY IMBALANCED AS NUMBER OF DEFULTERS IS VERY LESS IN TOTAL POPULATION. DATA IMBALANCE RATIO DEFULTER : NON-DEFULTER = 8 : 92 = 2 : 23.



HOUSING INFORMATION OF APPLICANT

ALL OF THE FEATURES HAVE VERY HIGH (47-70%) MISSING DATA PERCENTAGE. HENCE ALL THESE FEATURES CAN BE DROPPED.

MOST OF THE APPLICANTS LIVE IN HOUSE/APARTMENT

APPLICANTS LIVING WITH THEIR PARENTS OR IN RENTED APARTMENT HAVE HIGHER RATE OF DEFAULT.

		category	percentage
46	EMERGENCYSTATE_MODE	47.398304	
44	TOTALAREA_MODE	48.268517	
2	YEARS_BEGINEXPLUATATION_AVG	48.781019	
30	YEARS_BEGINEXPLUATATION_MEDI	48.781019	
16	YEARS_BEGINEXPLUATATION_MODE	48.781019	
35	FLOORSMAX_MEDI	49.780822	
7	FLOORSMAX_AVG	49.780822	
21	FLOORSMAX_MODE	49.780822	
43	HOUSETYPE_MODE	50.176091	
39	LIVINGAREA_MEDI	50.193326	
11	LIVINGAREA_AVG	50.193326	
25	LIVINGAREA_MODE	50.193326	
6	ENTRANCES_AVG	50.348768	
34	ENTRANCES_MEDI	50.348768	
20	ENTRANCES_MODE	50.348768	
28	APARTMENTS_MEDI	50.749729	
0	APARTMENTS_AVG	50.749729	
14	APARTMENTS_MODE	50.749729	
45	WALLSMATERIAL_MODE	50.840783	
19	ELEVATORS_MODE	53.295980	
5	ELEVATORS_AVG	53.295980	
33	ELEVATORS_MEDI	53.295980	
13	NONLIVINGAREA_AVG	55.179164	
41	NONLIVINGAREA_MEDI	55.179164	
27	NONLIVINGAREA_MODE	55.179164	
1	BASEMENTAREA_AVG	58.515956	
29	BASEMENTAREA_MEDI	58.515956	
15	BASEMENTAREA_MODE	58.515956	
23	LANDAREA_MODE	59.376738	
37	LANDAREA_MEDI	59.376738	
9	LANDAREA_AVG	59.376738	
3	YEARS_BUILD_AVG	66.497784	

	Value	Percentage of Defaulter
1	Rented apartment	12.313051
2	With parents	11.698113
3	Municipal apartment	8.539748
5	Co-op apartment	7.932264
0	House / apartment	7.795711
4	Office apartment	6.572411

ASSET DETAILS

MOST OF THE APPLICANT'S OWN REALTY

MOST OF THE APPLICANTS DO NOT OWN CARS

PEOPLE NOT OWNING REALITY AND CAR AND HAVE A SLIGHTLY HIGHER DEFAULT RATE THAN THE PEOPLE WHO OWN REALITY AND CAR.

DEFALTER OR NOT, MOST APPLICANTS HAVE CAR AGE BETWEEN 0-25 YEARS.

SINCE FOR BOTH TARGET VALUE, TREND IS SIMILAR, THIS FEATURE CAN DROP.

FAMILY RELATED INFO

DEFAULT RATE IS HIGHEST FOR CIVIL MARRIAGE AND SINGLE APPLICANTS

MOST OF THE APPLICANTS ARE MARRIED (AND/OR) NO CHILDREN (AND/OR) 2 FAMILY MEMBERS.

APPLICANTS WITH RELATIVELY GREATER NUMBER OF CHILDREN (AND/OR) FAMILY MEMBERS HAVE HIGHER DEFAULT PERCENTAGE.

FOR SOME OF THE CASES WHERE COUNT CHILDREN/FAMILY MEMBERS IS HIGH, AND THE DEFAULT RATE IS VERY HIGH OR VERY LOW. THESE CASES CANNOT BE TAKEN AS A CONCLUSION AS NUMBER OF APPLICANTS HAVING A LARGE FAMILY IS VERY LOW.

EDUCATION AND OCCUPATION INFO

MOST OF THE APPLICANTS ARE WORKING.

APPLICANTS ON MATERNITY LEAVE AND UNEMPLOYED HAS HIGHEST PERCENTAGE OF DEFALTER

BUSINESSMAN HAVE LOWEST (0) PERCENTAGE OF DEFALTER HOWEVER APPLICANTS OF INCOME TYPE('UNEMPLOYED', 'STUDENT', 'BUSINESSMAN', 'MATERNITY LEAVE') ARE VERY FEW IN THE DATASET TO CONTRIBUTE IN THE ANALYSIS.

APPLICANTS HAVING "LOWER SECONDARY" EDUCATION HAVE HIGHEST PERCENTAGE OF DEFALTER.

LOW SKILLED LABORERS HAVE VERY HIGH RATE OF DEFALTERS IN COMPARISON TO OTHER OCCUPATIONS.

FEMALE APPLICANTS ARE MORE THAN MALE APPLICANTS
DEFALUTER PERCENTAGE IS HIGHER FOR MALE APPLICANTS
PEOPLE OF AGE 30 HAVE HIGHER DEFAULT RATE

		Value	Percentage of Defaulter
7	Maternity leave	40.000000	
4	Unemployed	36.363636	
0	Working	9.588472	
2	Commercial associate	7.484466	
1	State servant	5.754965	
3	Pensioner	5.386366	
5	Student	0.000000	
6	Businessman	0.000000	

DEFAULT CASES ARE LESS FOR APPLICANTS MORE THAN 40 YEARS OLD.
MOST OF THE APPLICANTS ARE UNACCOMPANIED WHILE APPLYING FOR LOAN
NUMBER CASH LOANS IS QUITE HIGHER THAN REVOLVING LOANS
ALL WEEKDAYS HAVE SIMILAR NUMBER OF APPLICANTS THAN WEEKEND (SATURDAY AND SUNDAY)

BOXPLOT IS SHOWING THE OUTLIERS FOR INCOME AND ANNUITY, THERE ARE FEW ENTRIES HAVING VERY LARGE ANNUITY AND INCOME THAN OTHERS.



TOP 10 CORRELATION OF DEFAULTERS

AMT_REQ_CREDIT_BUREAU_YEAR	AMT_REQ_CREDIT_BUREAU_YEAR	1.000000
OBS_60_CNT_SOCIAL_CIRCLE	OBS_30_CNT_SOCIAL_CIRCLE	0.998269
BASEMENTAREA_AVG	BASEMENTAREA_MEDI	0.998250
COMMONAREA_AVG	COMMONAREA_MEDI	0.998107
YEARS_BUILD_MEDI	YEARS_BUILD_AVG	0.998100
NONLIVINGAPARTMENTS_AVG	NONLIVINGAPARTMENTS_MEDI	0.998075
FLOORSMIN_AVG	FLOORSMIN_MEDI	0.997825
LIVINGAPARTMENTS_AVG	LIVINGAPARTMENTS_MEDI	0.997668
FLOORSMAX_MEDI	FLOORSMAX_AVG	0.997187
NONLIVINGAPARTMENTS_MEDI	NONLIVINGAPARTMENTS_MODE	0.997032
ENTRANCES_MEDI	ENTRANCES_AVG	0.996700

TOP 10 CORRELATION OF NON DEFAULTERS

AMT_REQ_CREDIT_BUREAU_YEAR	AMT_REQ_CREDIT_BUREAU_YEAR	1.000000
YEARS_BUILD_AVG	YEARS_BUILD_MEDI	0.998522
OBS_60_CNT_SOCIAL_CIRCLE	OBS_30_CNT_SOCIAL_CIRCLE	0.998508
FLOORSMIN_MEDI	FLOORSMIN_AVG	0.997202
FLOORSMAX_MEDI	FLOORSMAX_AVG	0.997018
ENTRANCES_MEDI	ENTRANCES_AVG	0.996899
ELEVATORS_AVG	ELEVATORS_MEDI	0.996161
COMMONAREA_MEDI	COMMONAREA_AVG	0.995857
LIVINGAREA_AVG	LIVINGAREA_MEDI	0.995568
APARTMENTS_AVG	APARTMENTS_MEDI	0.995163
BASEMENTAREA_MEDI	BASEMENTAREA_AVG	0.994081

RESULT - THIS DATA IS HIGHLY IMBALANCED AS NUMBER OF DEFAULTERS IS VERY LESS IN TOTAL POPULATION.

MOST OF THE APPLICANTS LIVE IN HOUSE/APARTMENT II. APPLICANTS LIVING WITH THEIR PARENTS OR IN RENTED APARTMENT HAVE HIGHER RATE OF DEFAULT.

SOCIAL CIRCLE INFO: THE FEATURES SHOW SIMILAR TREND FOR DEFAULTERS AND NON-DEFAULTERS, CAN BE DROPPED.

ASSET INFO - MOST OF THE APPLICANT'S OWN REALTY II. MOST OF THE APPLICANTS DO NOT OWN CARS III. PEOPLE NOT OWNING REALITY AND CAR AND HAVE A SLIGHTLY HIGHER DEFAULT RATE THAN THE PEOPLE WHO OWN REALITY AND CAR.

XYZ ADS AIRING REPORT ANALYSIS

DESCRIPTION:

IN SIMPLER TERMS, ADVERTISING ANALYTICS CAN BE REFERRED TO AS THE USE OF ANALYTICAL DATA AND TOOLS THAT HELP BUSINESSES AND MARKETERS EFFICIENTLY MONITOR THEIR OMNICHANNEL MARKETING EFFORTS. AS THESE DATA SETS OFFER ACTIONABLE INSIGHTS, MARKETERS CAN USE THEM TO REASSURE THAT THE CAMPAIGNS THEY RUN ARE TARGETED TO THE RIGHT AUDIENCE AND USE THE RIGHT MEDIUM TO DO SO.

AS PREVIOUSLY MENTIONED, ORGANIZATIONS NEVER RESORT TO A SINGLE MARKETING CHANNEL AS THEIR TARGET AUDIENCES WILL BE SPREAD ACROSS MULTIPLE CHANNELS. TO BE MORE PRECISE, ORGANIZATIONS DON'T HAVE A CHOICE OTHER THAN TO RUN CROSS-CHANNEL CAMPAIGNS AS MULTI-CHANNEL CUSTOMERS TEND TO SPEND TWO TO FOUR TIMES MORE THAN SINGLE-CHANNEL CUSTOMERS.

RESULTS

A. WHAT IS POD POSITION? DOES THE POD POSITION NUMBER AFFECT THE AMOUNT SPENT ON ADS FOR A SPECIFIC PERIOD OF TIME BY A COMPANY? (EXPLAIN IN DETAILS WITH EXAMPLES FROM THE DATASET PROVIDED)

AN AD POD IS A GROUP OF ADS THAT ARE SEQUENCED TOGETHER TO BE PLAYED BACK-TO-BACK WITHIN A SINGLE AD BREAK/PLACEMENT, SIMILAR TO AD BREAKS IN TRADITIONAL LINEAR TV. AD PODS GIVE PUBLISHERS THE OPPORTUNITY TO MAXIMISE REVENUE FROM EACH AD BREAK AND GIVE ADVERTISERS MORE CONTROL OVER AD POSITIONING.

THEY ALLOW PUBLISHERS TO RETURN MULTIPLE ADS FROM A SINGLE AD REQUEST, AND THEN THOSE ADS ARE PLAYED IN SEQUENCE.
FOUR REASONS WHY ADVERTISERS AND PUBLISHERS USE AD PODS

- 1) THEY OFFER MORE CONTROL AD PODS HELP ADVERTISERS AVOID RUNNING ADS ALONGSIDE DIRECT COMPETITORS, ENSURING THAT THEIR OFFERING STANDS OUT TO VIEWERS AND THAT THEIR MESSAGE DOESN'T GET SATURATED
- 2) AD PODS OFFER A BETTER WAY TO MONETIZE LONG-FORM CONTENT PUBLISHERS WITH LONGER-FORM CONTENT CAN LEVERAGE THE CONTROLS OFFERED BY AD PODDING TO SET UP MORE ADVANCED MONETIZATION STRATEGIES FOR THEIR STREAMING CONTENT.

3)- AD PODS ALLOW PUBLISHERS TO MEET BUYERS' NEEDS WITHOUT AN AD POD IN PLACE, ADVERTISERS HAVE HISTORICALLY BEEN SPINNING A WHEEL OF CHANCE WHEN THEY BOUGHT INTO STREAMING APPS ON CTV. THEY HAVE HAD VERY LITTLE CONTROL OVER FREQUENCY OR THE POSITION OF THEIR AD WITHIN THE AD BREAK.

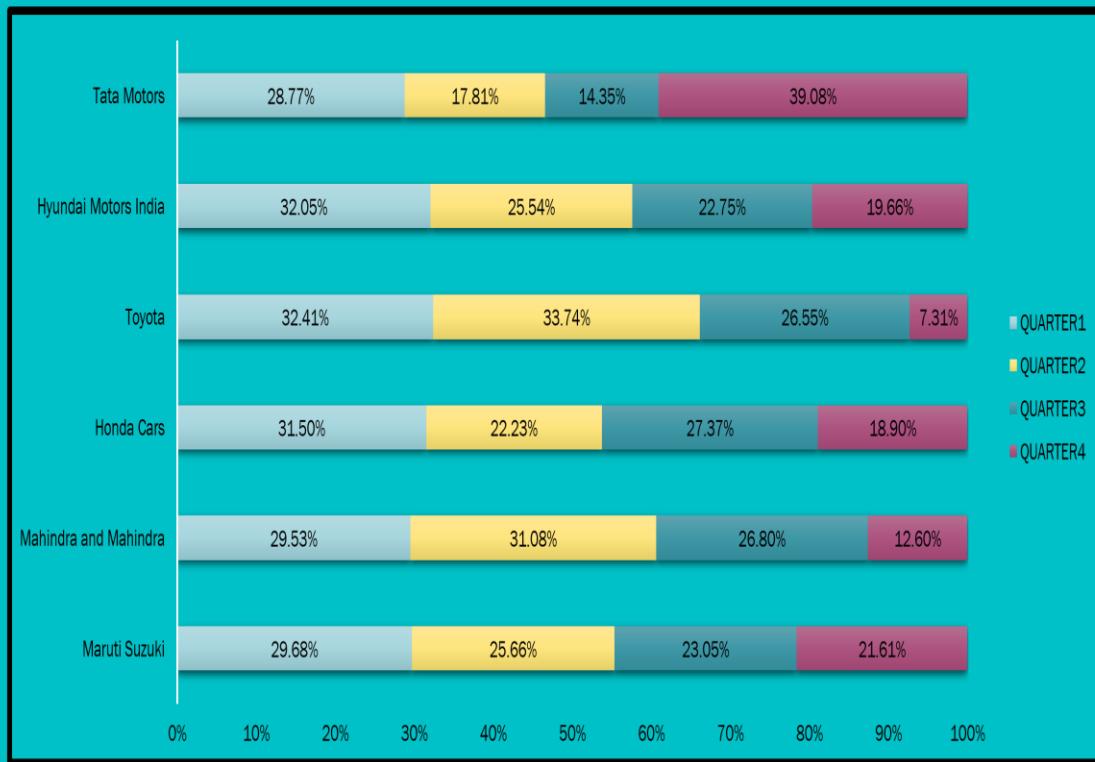
4)-AD PODS IMPROVE THE USER EXPERIENCE AD PODS ENABLE A BETTER USER EXPERIENCE BECAUSE THE ADS ARE TAILED TO VIEWERS' PREFERENCES.

Pod Positions	Honda Cars	Hyundai Motors India	Mahindra and Mahindra	Maruti Suzuki	Tata Motors	Toyota
1	26%	28%	26%	37%	30%	26%
2	23%	24%	21%	21%	22%	22%
3	20%	20%	21%	18%	20%	21%
4	17%	16%	18%	14%	16%	18%
5	14%	12%	15%	11%	13%	14%

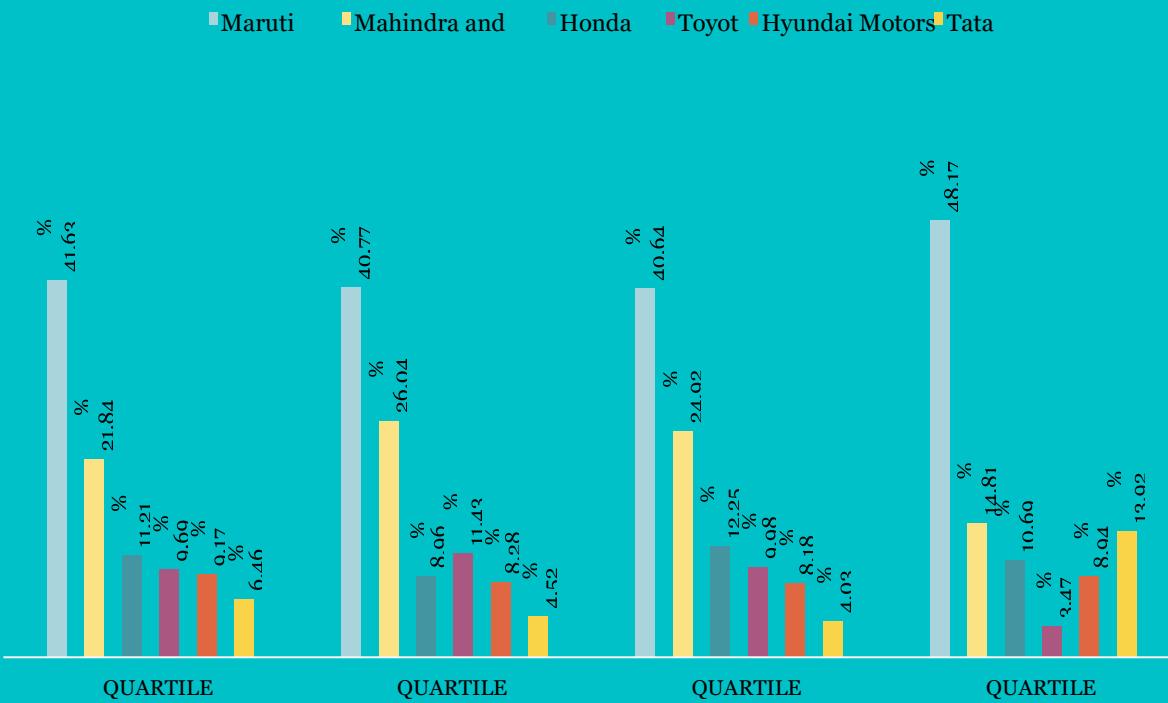
- THE ABOVE ANALYSIS SHOWS THAT ALL BRANDS HAVE THE HIGHEST CONTRIBUTION OF EQ UNITS IN 1 ST POD_POSITION AMONG THE TOP 5 POD_POSITIONS.



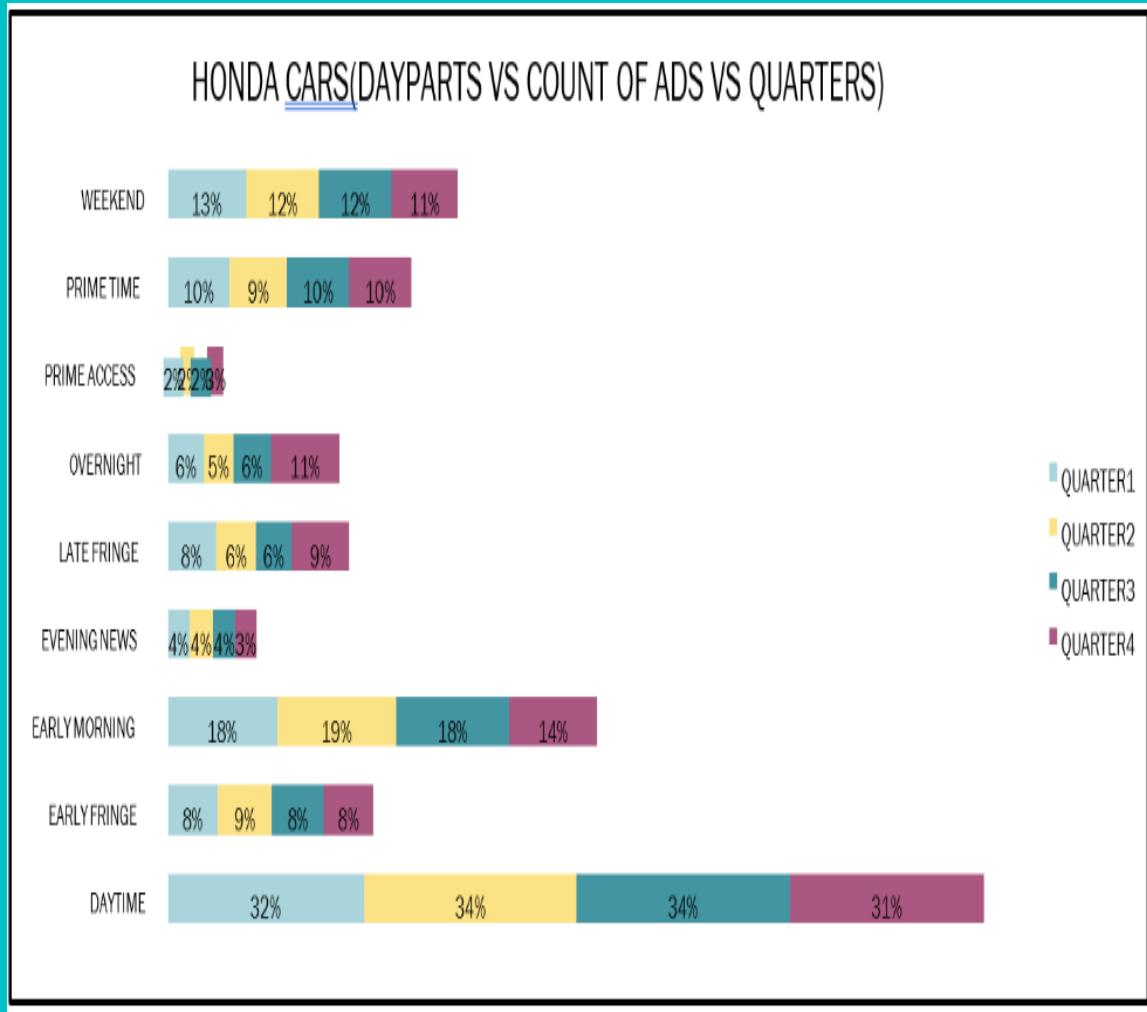
DISTRIBUTION OF ADS AS PER SUM OF EQ UNITS



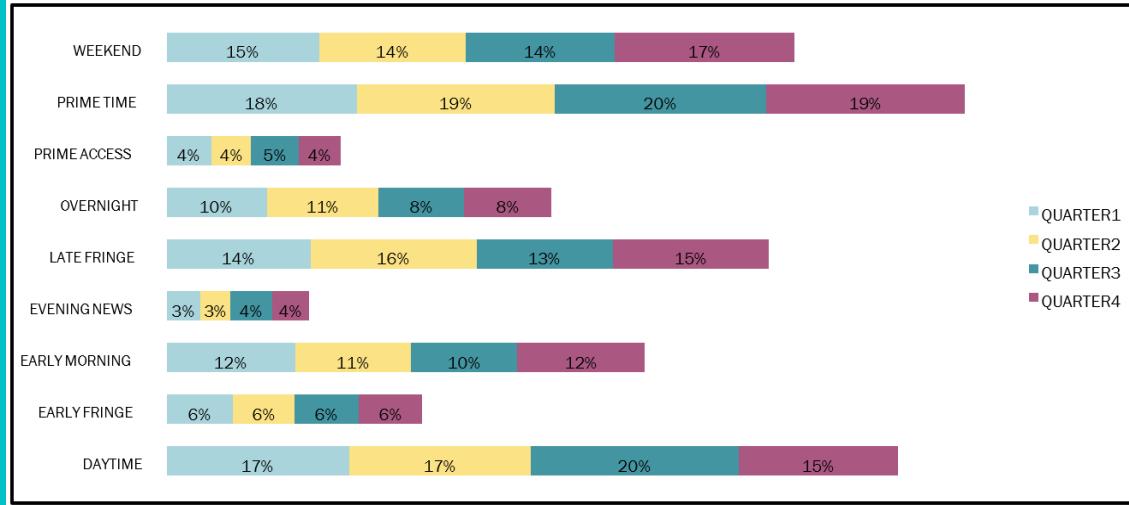
DISTRIBUTION OF ADS AS PER THEIR DURATION



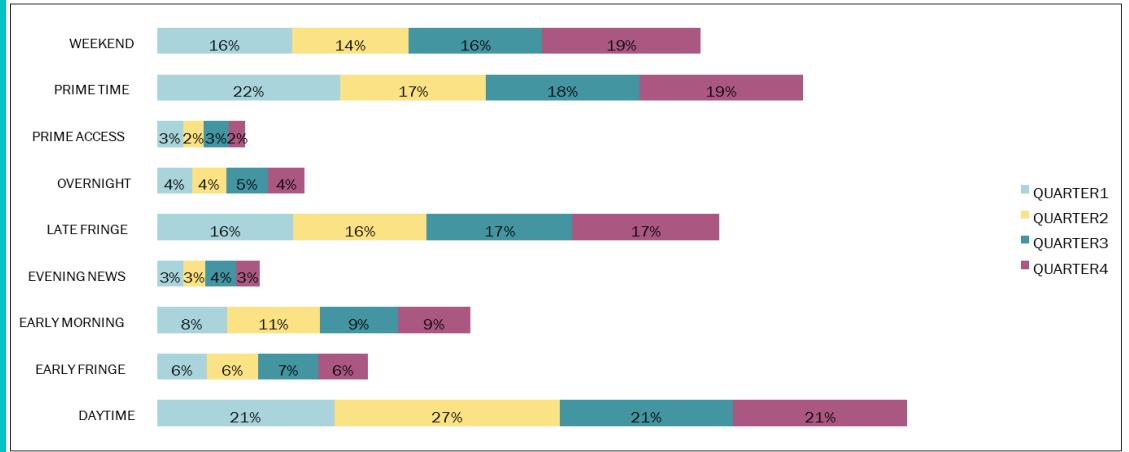
CONDUCT A COMPETITIVE ANALYSIS FOR THE BRANDS AND DEFINE ADVERTISEMENT STRATEGY OF DIFFERENT BRANDS AND HOW IT DIFFERS ACROSS THE BRANDS.



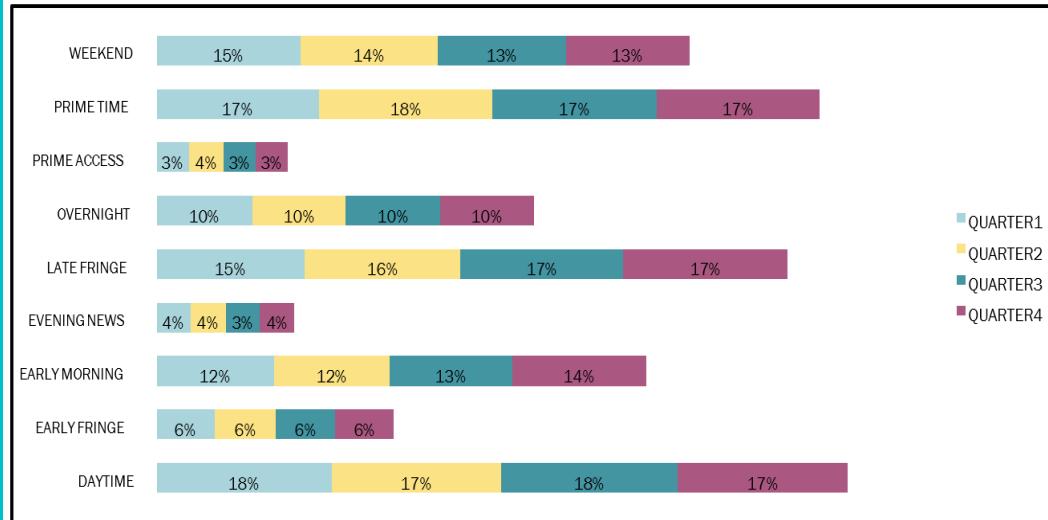
HYUNDAI MOTORS INDIA(DAY PARTS VS COUNT OF ADS VS QUARTERS)



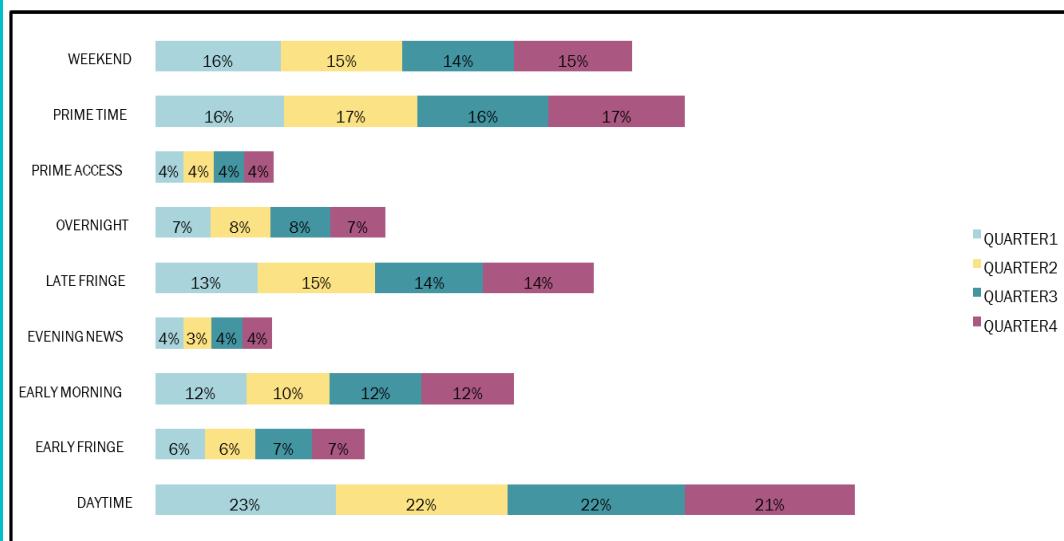
MAHINDRA AND MAHINDRA(DAY PARTS VS COUNT OF ADS VS QUARTERS)



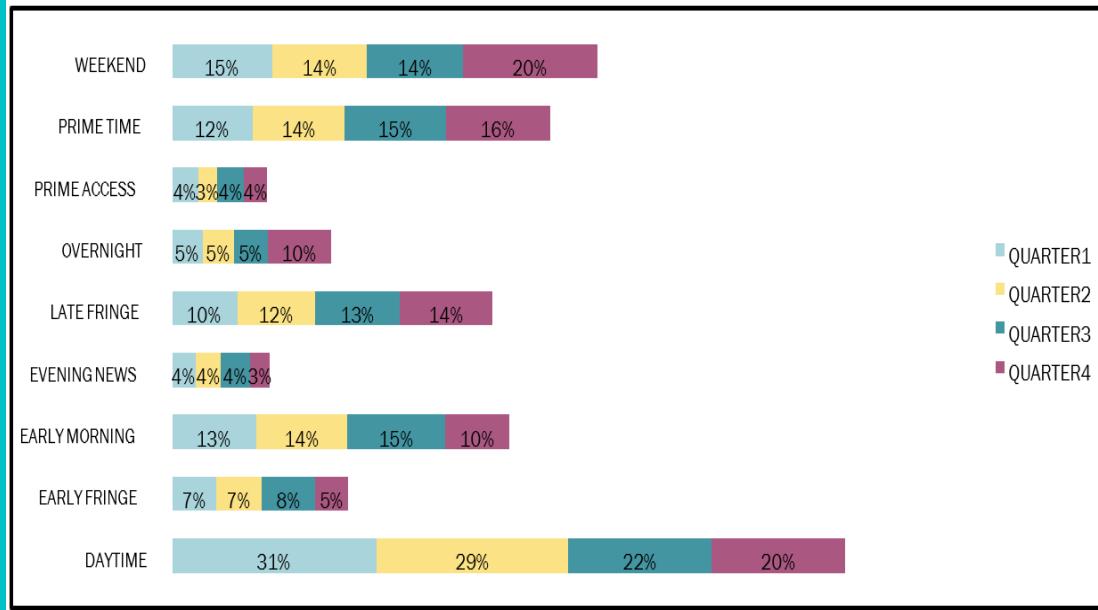
MARUTI SUZUKI (DAY PARTS VS COUNT OF ADS VS QUARTERS)



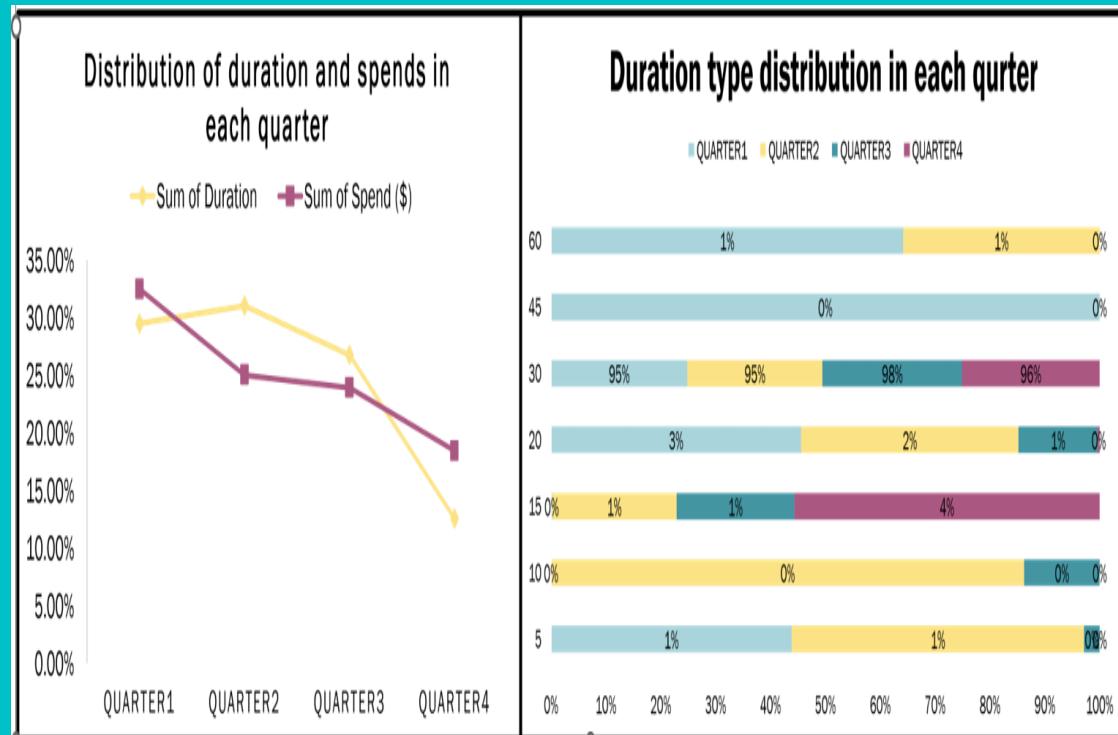
TATA MOTORS (DAYPARTS VS COUNT OF ADS VS QUARTERS)

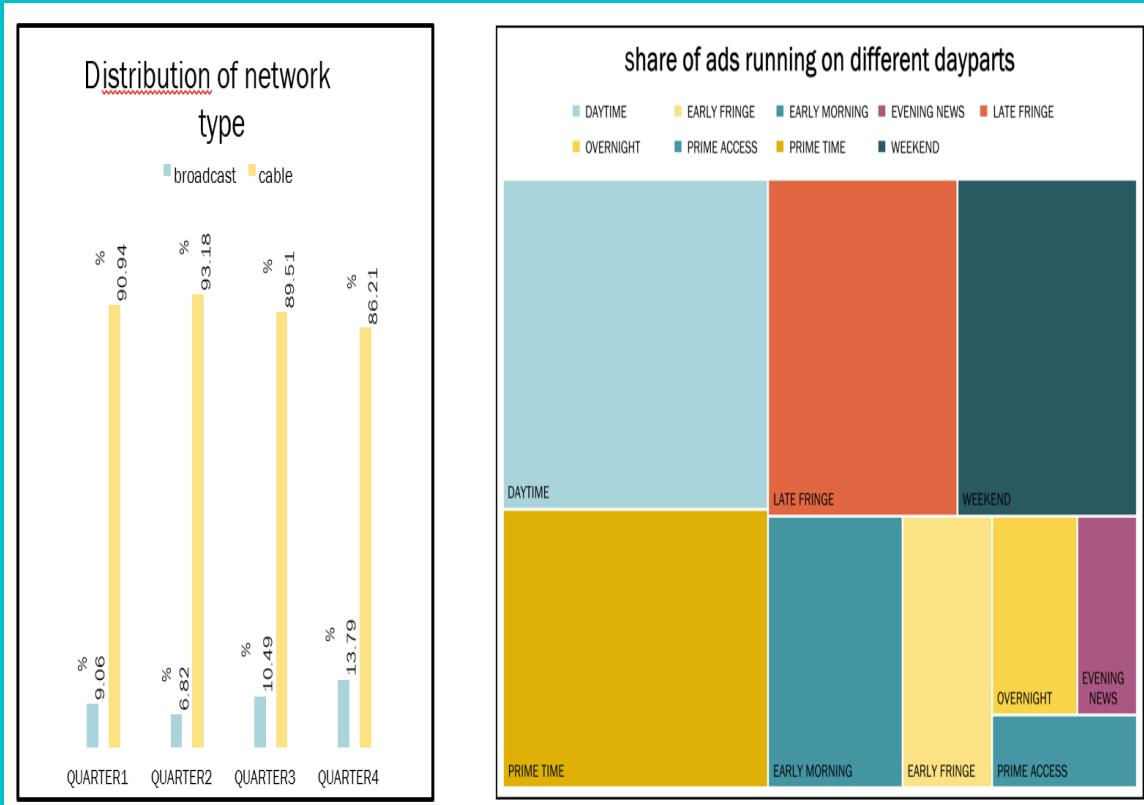


TOYOTA (DAYPARTS VS COUNT OF ADS VS QUARTERS)

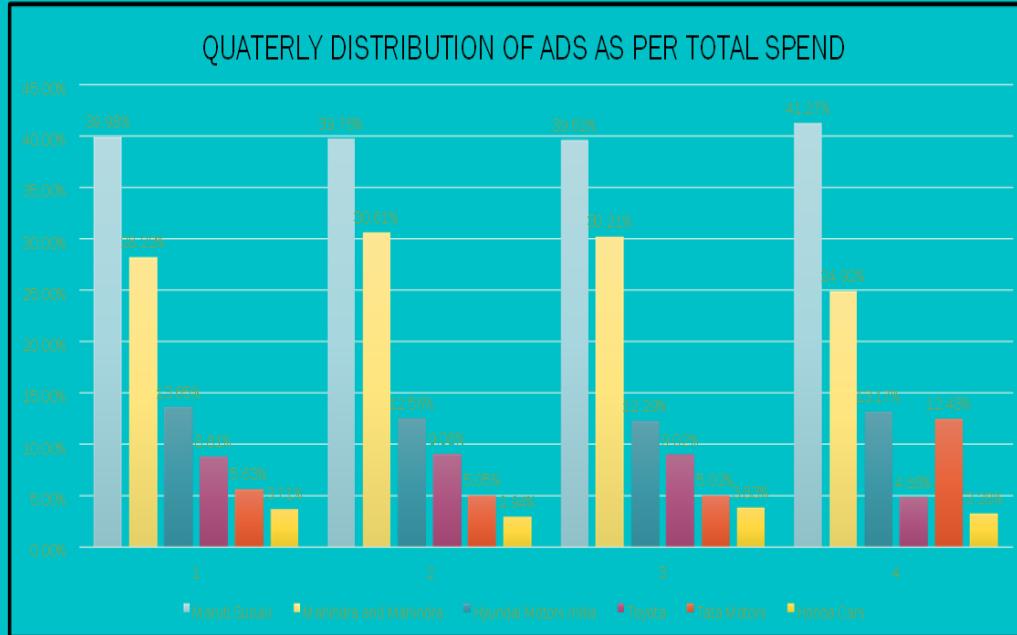


MAHINDRA AND MAHINDRA WANTS TO RUN A DIGITAL AD CAMPAIGN TO COMPLEMENT ITS EXISTING TV ADS IN Q1 OF 2022. BASED ON THE DATA FROM 2021, SUGGEST A MEDIA PLAN TO THE CMO OF MAHINDRA AND MAHINDRA. WHICH AUDIENCE SHOULD THEY TARGET? *ASSUME XYZ ADS HAS THE AD VIEWERSHIP DATA AND TV VIEWERSHIP FOR THE PEOPLE IN INDIA.





WHAT IS THE SHARE OF VARIOUS BRANDS IN TV AIRINGS AND HOW HAS IT CHANGED FROM Q1 TO Q4 IN 2021?



2022 Q1 MEDIA PLAN FOR

MARKETING STRATEGIES

- **MAHINDRA AND MAHINDRA SHOULD TARGET AUDIENCE FROM NORTHEAST AND SOUTHERN PARTS OF INDIA AS A TARGET CUSTOMER IN Q1 2022 TO ENHANCE THEIR BUSINESS.**
- **SINCE THE MAIN PRODUCT OF MAHINDRA AND MAHINDRA IS PERSONAL AUTOMOBILES AND INSURANCE SO THEY NEED TO INCREASE THEIR MEDIA SPENT AND ADS DURATION ON SPORTS CHANNEL, NEWS CHANNEL AND CORPORATES ENTERTAINMENT NETWORK AS COMPARED TO FAMILY ENTERTAINMENT AND CHILD ENTERTAINMENT IN ORDER TO MAXIMIZE THEIR REACH ON TARGET AUDIENCE.**

Distribution of amount spent by different brands on a monthly basis

Months	Honda Cars	Hyundai Motors India	Mahindra and Mahindra	Maruti Suzuki	Tata Motors	Toyota
January	10.35%	11.66%	11.34%	12.47%	9.79%	13.75%
February	13.85%	12.44%	9.30%	9.62%	10.90%	10.80%
March	10.95%	10.45%	11.89%	10.67%	6.52%	11.26%
April	5.04%	7.49%	8.03%	7.55%	5.30%	8.44%
May	6.85%	7.44%	8.15%	7.29%	6.55%	8.69%
June	7.89%	7.65%	8.86%	8.29%	5.46%	9.07%
July	10.15%	7.09%	8.25%	8.17%	6.56%	8.64%
August	8.06%	7.30%	7.91%	7.26%	2.58%	9.79%
September	6.81%	7.02%	7.80%	6.91%	7.56%	6.80%
October	6.20%	7.27%	6.71%	8.00%	11.56%	4.10%
November	5.84%	6.67%	5.45%	6.35%	13.27%	3.83%
December	8.00%	7.51%	6.31%	7.41%	13.95%	4.82%

ABC CALL VOLUME TREND ANALYSIS

PROJECT DESCRIPTION:

ABC IS A CALL CENTRE WHICH HAS A CUSTOMER EXPERIENCE TEAM FOR THE VOICE PROCESS. TYPICALLY, THESE TEAMS FULFIL VARIOUS ROLES AND RESPONSIBILITIES SUCH AS: CUSTOMER EXPERIENCE PROGRAMS (CX PROGRAMS), DIGITAL CUSTOMER EXPERIENCE, DESIGN AND PROCESSES, INTERNAL COMMUNICATIONS, VOICE OF THE CUSTOMER (VOC), USER EXPERIENCES, CUSTOMER EXPERIENCE MANAGEMENT, JOURNEY MAPPING, NURTURING CUSTOMER INTERACTIONS, CUSTOMER SUCCESS, CUSTOMER SUPPORT, HANDLING CUSTOMER DATA, LEARNING ABOUT THE CUSTOMER JOURNEY. I HAVE BEEN PROVIDED WITH THE DATA OF ABC CALL CENTRE FOR THE LAST 23 DAYS AND I SHOULD ANALYZE THE DATA AND HELP THE COMPANY ANSWER SOME OF THE BUSINESS QUESTIONS.

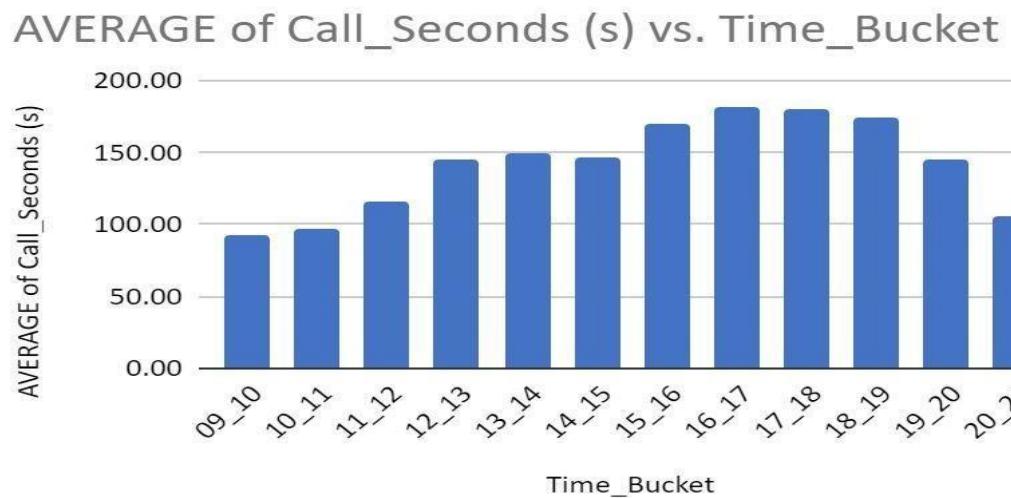
ABOUT THE DATASET :

- THERE ARE 117989 ROWS AND 13 COLUMNS.
- AGENT NAME IS A CATEGORICAL COLUMN CONTAINING NAMES OF THE AGENTS ATTENDING THE PHONE CALL.
- AGENT ID COLUMN CONTAINS THE AGENT ID AND THEY ALSO CONTAIN A LOT OF NULL VALUES.
- CUSTOMER COLUMN CONTAINS THE CONTACT NUMBER OF THE CUSTOMER.
- QUEUE TIME CONTAINS THE SECONDS CUSTOMER HAS WAITED BEFORE THE AGENT HAS PICKED THE CALL.
- DATE & TIME, TIME & TIME_BUCKET ARE USED FOR TIME INTELLIGENCE.
- DURATION AND CALL_SECONDS CONTAINS THE DURATION OF THE CALL.
- CALL STATUS AND WRAPPED ARE ALSO CATEGORICAL COLUMNS.

CLEANING THE DATASET :

- THE AGENT NAME AND AGENT ID HAS NULL VALUES.
- THE COLUMNS RINGING AND IVR DURATION ARE NOT USED SINCE RINGING HAS ONLY ONE VARIABLE AND IVR DURATION IS NOT OF ANY USE TO OUR ANALYSIS.

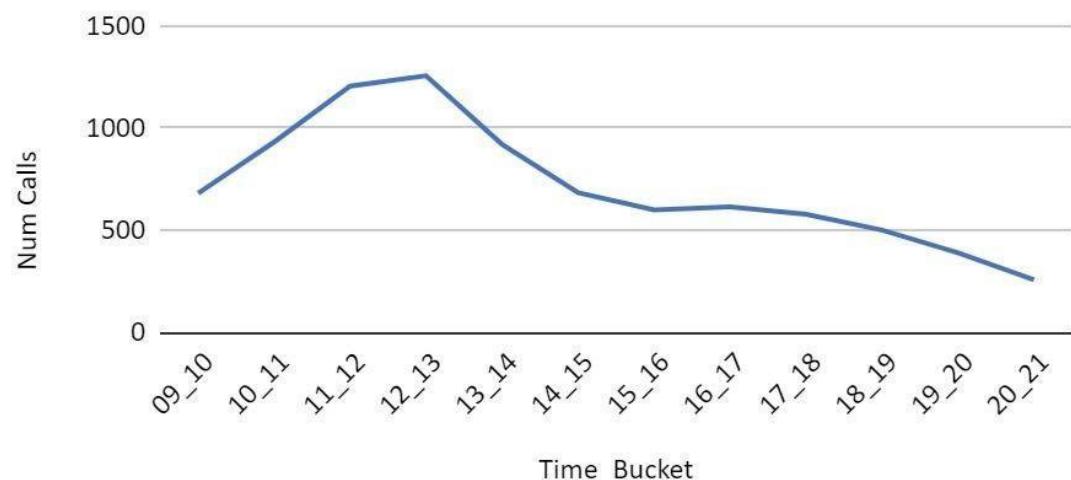
AVERAGE CALL TIME DURATION FOR EACH TIME BUCKET



IT CAN BE SEEN THAT THE AVERAGE CALL DURATION INCREASES FROM 9 TO 5 AND THEN DECREASES.

NO. OF CALLS FOR EACH TIME BUCKET

Num Calls vs. Time_Bucket



IT CAN BE SEEN THAT THE NUMBER OF CALLS IN A DAY INCREASES FROM 9 TO 12 AND THEN IT STARTS DECREASING.

DECREASE THE DROP RATE

EXECUTIVE NAME	DROPPED	ANSWERED	TRANSFERRED
#N/A	34198		
GRAND TOTAL	34403	82452	1133

IT CAN BE SEEN THAT OUT OF 34403 CALLS DROPPED, 34198 CALLS HAVE NOT BEEN ANSWERED BY ANYONE. THIS CAN ALSO MEAN THAT THE AGENTS WERE BUSY IN ANY OTHER CALLS AND SO, THEY WERE UNABLE TO TEND TO THESE CALLS. THIS MEANS THAT WE NEED TO INCREASE THE MAN POWER TO REDUCE THE DROP RATE.

TIME_BUCK ET	NUMBER OF AGENTS	ABANDO N RATE	CALLS EACH DAY	AVG CALL DURATION	AVG QUEUE TIME
09_10	42	53.70%	416.87	92.01	82.86
10_11	51	51.91%	578.83	97.42	83.25
11_12	59	41.21%	635.91	116.78	72.32
12_13	60	24.29%	550.09	144.73	41.66
13_14	58	22.64%	502.65	149.54	41.80
14_15	60	23.44%	459.17	146.97	43.60
15_16	58	13.25%	398.22	169.90	29.88
16_17	58	8.50%	382.09	181.44	23.54
17_18	58	9.18%	371.04	179.72	23.75
18_19	59	12.89%	314.70	174.32	34.09
19_20	52	28.59%	281.00	144.58	58.69
20_21	27	47.68%	239.35	105.95	75.28
GRAND TOTAL	66	29.16%	5,129.91	139.53	52.17

IT CAN BE INFERRED THAT THE ABANDON RATES ARE HIGH ON PARTICULAR TIMES OF A DAY AND LOWEST ON A PARTICULAR TIME IN A DAY. EVEN THOUGH THE CALLS ARE LOW DURING THE 9_10 BUCKET, THE DROP RATE IS VERY HIGH. WE HAVE THE AVERAGE CALL DURATION IN EACH TIME BUCKET. USING THAT, WE CAN CALCULATE APPROXIMATELY HOW MUCH TIME THE CALL CENTRE AGENTS ARE SPENDING TALKING TO THE CUSTOMERS IN TOTAL IN A PARTICULAR TIME BUCKET. BY USING THAT DATA, WE CAN ARRIVE AT THE APPROXIMATE NUMBER OF AGENTS REQUIRED TO HAVE A DROP RATE LESSER THAN 10 %. IT CAN ALSO BE SEEN THAT THE TIME BUCKETS IN WHICH THE AVG CALL DURATION IS HIGH AND AVG QUEUE TIME IS LOW HAVE LESSER DROP RATES.

NOTE: A TOLERANCE LEVEL NEEDS TO BE ADDED IN EVERY CALCULATION TO ACCOUNT FOR THE ERRORS.

TIME_BUCKET	9 SHIFT	10 SHIFT	12 SHIFT	AGENTS REQ	AGENTS REQ CALCULATED
09_10	50	0	0	50	44.88
10_11	50	20	0	70	64.04
11_12	50	20	0	70	73.96
12_13	25	10	40	75	62.65
13_14	0	20	40	60	58.47
14_15	50	0	20	70	53.57
15_16	25	20	40	85	48.84
16_17	50	20	0	70	47.60
17_18	50	10	20	80	46.01
18_19	0	20	30	50	39.48
19_20	0	0	40	40	34.78
20_21	0	0	30	30	26.58

AGENTS REQUIRED IS CALCULATED BY USING THE FORMULA

AGENTS REQ CALCULATED = 1.1 * CALLS EACH DAY * 2 * (AVG CALL DURATION + AVG QUEUE TIME) WHERE 1.1 AND 2 ARE TOLERANCE OF NO OF CALLS AND WAITING TIME RESPECTIVELY.

SO THE TOTAL MAN POWER REQUIRED = 9 AM SHIFT + 10 AM SHIFT + 12 PM SHIFT = 50 + 20 + 40 = 110

NEW MANPOWER TO BE ADDED = REQ MAN POWER - AVAILABLE EMPLOYEES = 110 - 66 = 44

IT IS ADVISED TO ADD 44 NEW EMPLOYEES TO REDUCE THE ABANDON RATE FROM 30 % TO 10 %

4. ORGANIZE NIGHT CALLS

MEN REQUIRED CAN BE CALCULATED FROM THE SAME FORMULA GIVEN IN THE PREVIOUS SLIDE. FOR THE NIGHT CALLS TOO, MEN REQUIRED CAN BE FOUND FROM THE SAME FORMULA. CALLS/ DAY CAN BE FOUND FROM THE DATA FOR THE WHOLE 24 HOURS. A SHIFT PLAN HAS BEEN DRAFTED IN WHICH THE REQUIRED NUMBER OF AGENTS IN THAT SPECIFIC PERIOD OF TIME HAS BEEN USED TO CAREFULLY DECIDED BASED ON THE EXPECTED NUMBER OF CALLS ALONG WITH SOME TOLERANCE.

8 AM SHIFT

TIME BUCKET	MEN REQUIRED	TOTAL MEN	8 AM SHIFT	5 PM SHIFT	2 AM SHIFT
08_09	42.93	80	60	0	20
09_10	44.88	90	70	0	20
10_11	64.04	70	70	0	0
11_12	73.96	70	70	0	0
12_13	62.65	60	60	0	0
13_14	58.47	60	60	0	0
14_15	53.57	60	60	0	0
15_16	48.84	50	50	0	0
16_17	47.60	50	50	0	0

5 PM SHIFT

TIME BUCKET	MEN REQUIRED	TOTAL MEN	8 AM SHIFT	5 PM SHIFT	2 AM SHIFT
17_18	46.01	50	0	50	0
18_19	39.48	50	0	50	0

19_20	34.78	50	0	50	0
20_21	26.58	50	0	50	0
21_22	25.76	25	0	25	0
22_23	25.76	25	0	25	0
22-00	17.17	50	0	50	0
00_01	17.17	50	0	50	0
01_02	8.59	70	0	50	20

2 AM SHIFT

TIME BUCKET	MEN REQUIRED	TOTAL MEN	8 AM SHIFT	5 PM SHIFT	2 AM SHIFT
02_03	8.59	70	0	50	20
03_04	8.59	20	0	0	20
04_05	8.59	10	0	0	10
05_06	8.59	10	0	0	10
06_07	25.76	20	0	0	20
07_08	34.34	20	0	0	20
08_09	34.34	20	0	0	20
02_03	42.93	80	60	0	20
09_10	44.88	90	70	0	20

SO, THE TOTAL NUMBER OF MAN POWER NEEDED IS 8 AM SHIFT + 5 PM SHIFT + 2 AM SHIFT = $70 + 50 + 20 = 140$.

AFTER HIRING $140 - 66 = 74$ AGENTS, WE CAN EXPECT THE CALL ABANDON RATES TO LESSER THAN 10 %

CONCLUSION

IN CONCLUSION, I WOULD LIKE TO TELL THAT AFTER DOING A THOROUGH ANALYSIS WE WERE ABLE TO DERIVE THE INSIGHTS FROM THE DATA AND WAS ABLE TO PLOT VARIOUS GRAPHS USING THAT DATA.

THE DATA WHICH ONCE LOOKED USELESS GAVE SOME VERY USEFUL INSIGHTS.