



# IMDB Movie Analysis

- Abhishek Shukla

# Project Description

The IMDB Movie Analysis Project is based on the world wide realised movie records, it is known also as an analytical website on this website we collect information about movies from real world and Analyze them on the basis of users reviews, likes and other a lot of parameter so the project begins with some csv dataset which is totally row in this project we are going to clean the data set and Analyze the records given task by stakeholders to provide some meaningful insights.

# Approach of the Analysis

- ✓ First Clean the dataset to start the process of Analysis.
- ✓ Review the dataset and understand the different columns and their meanings.
- ✓ Use the provided guiding questions to frame the problem you want to shed light on.
- ✓ Compile the findings from the analysis into a detailed report.
- ✓ Create Visuals to better understanding of the stakeholders.
- ✓ Present the data story in a clear and concise manner, making it easy for stakeholders to understand.



# Tech-Stack Used In The Analysis

- ✓ Excel 365 for complete analysis and visualization.
- ✓ Power Point Presentation 365 is used to create reports with figures.

# Task A : Cleaning the data

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.)

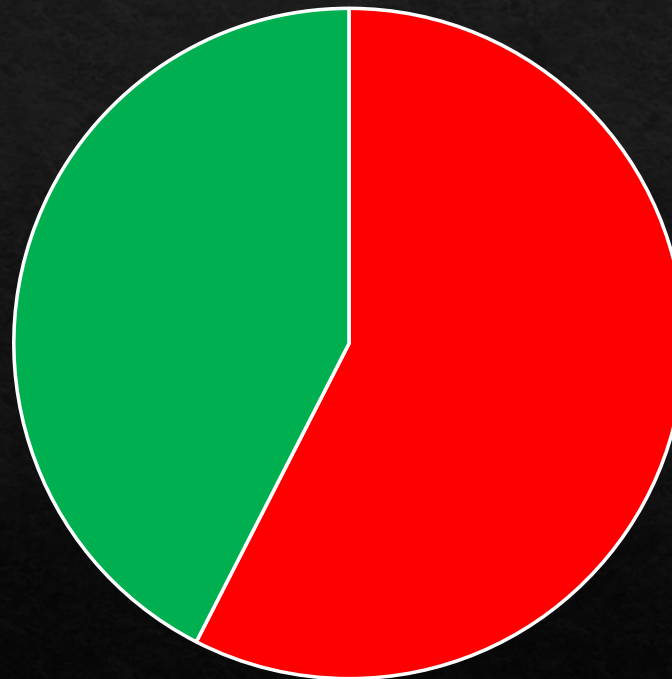
Your task: Clean the data

# Task A - Data Cleaning

## Data Cleaning Result

Now Starts data Analysis With Clean Data to extract some meaningful insights.

Coloumn after data cleaning, 3709, 42%



In row dataset I found lots of null values and duplicates after removing this I get 3709 out of 5025 records, which will serve as the starting point for analysis.

Row table Coloumn, 5025, 58%

■ Row table Coloumn    ■ Coloumn after data cleaning

# Task B : Movies with highest profit

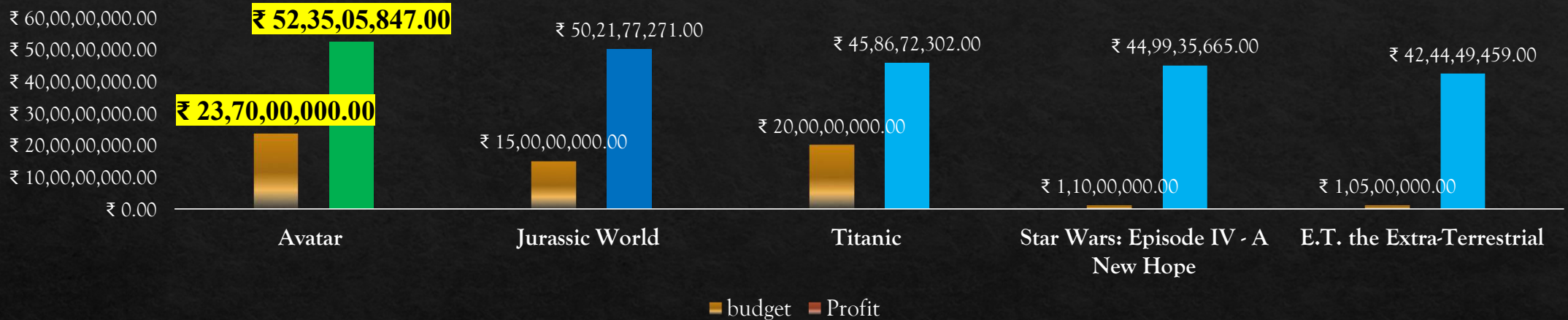
Movies with highest profit : Create a new column called profit which contains the difference of the two columns: gross and budget. Sort the column using the profit column as reference. Plot profit (y-axis) vs budget (x- axis) and observe the outliers using the appropriate chart type.

Your task: Find the movies with the highest profit?



# Task B - Find the movies with the highest profit?

## TOP 5 MOST PROFITABLE MOVIES



Insights :

These are the top 5 profitable movies in given database and the most profitable movie is "Avatar".

The screenshot shows an Excel spreadsheet with a formula bar at the top. The formula bar contains the formula `=[@gross]-[@budget]`, which is circled in red. The spreadsheet has columns for various movie attributes, including country, content\_rating, budget, title\_year, actor\_2\_facebook\_likes, imdb\_score, aspect\_ratio, movie\_facebook\_likes, and Profit. The Profit column is highlighted in green.

country	content_rating	budget	title_year	actor_2_facebook_likes	imdb_score	aspect_ratio	movie_facebook_likes	Profit
USA	PG-13	237000000	2009	936	7.9	1.78	33000	523505847
USA	PG-13	150000000	2015	2000	7	2	150000	502177271
USA	PG-13	200000000	1997	14000	7.7	2.35	26000	458672302
USA	PG	110000000	1977	1000	8.7	2.35	33000	449935665
USA	PG	105000000	1982	725	7.9	1.85	34000	424449459
USA	PG-13	220000000	2012	21000	8.1	1.85	123000	403279547
USA	G	45000000	1994	886	8.5	1.66	17000	377783777
USA	PG	115000000	1999	14000	6.5	2.35	13000	359544677
USA	PG-13	185000000	2008	13000	9	2.35	37000	348316061
USA	PG-13	78000000	2012	14000	7.3	2.35	140000	329999255
USA	R	58000000	2016	805	8.1	2.35	117000	305024263
USA	PG-13	130000000	2013	14000	7.6	2.35	82000	294645577
USA	PG-13	63000000	1993	610	8.1	1.85	19000	293784000
USA	PG	76000000	2013	2000	7.5	1.85	56000	292049635
USA	R	58800000	2014	962	7.3	2.35	112000	291323553
USA	G	94000000	2003	939	8.2	1.85	11000	286838870
USA	PG	150000000	2004	309	7.2	1.85	0	286471036
USA	PG-13	94000000	2003	857	8.9	2.35	16000	283019252
USA	PG	32500000	1983	1000	8.4	2.35	14000	276625409

There is Formula How To Calculate Profit.



# Task C : Top 250

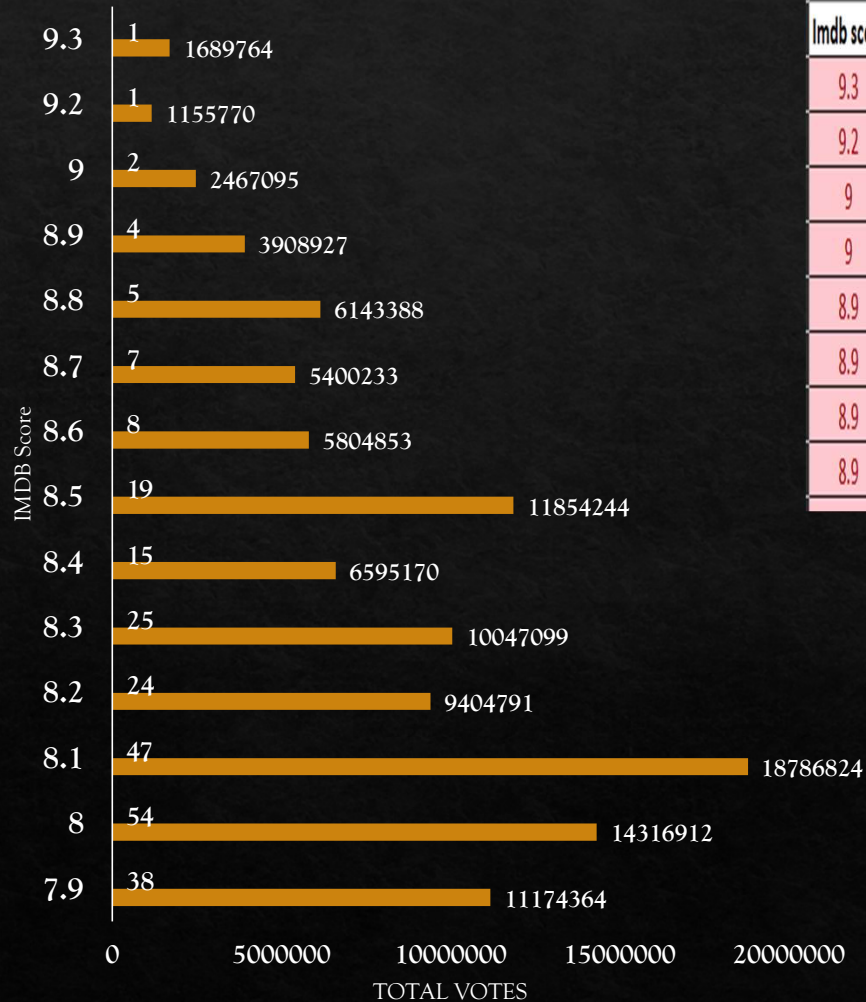
Create a new column IMDb\_Top\_250 and store the top 250 movies with the highest IMDb Rating (corresponding to the column : imdb\_score). Also make sure that for all of these movies, the num\_voted\_users is greater than 25,000. Also add a Rank column containing the values 1 to 250 indicating the ranks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

# Task C - Find IMDB Top 250

Top 250 Movies

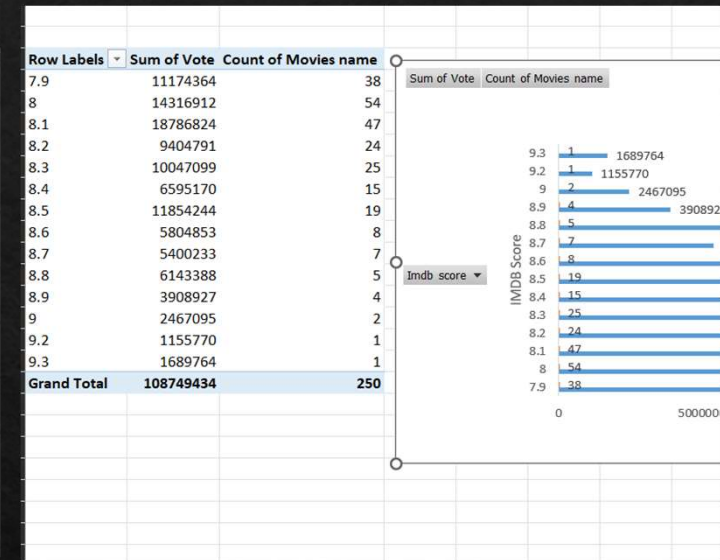


A	B	C	D	E
Imdb score	Top 250	Vote	Rank	
9.3	The Shawshank Redemption	1689764	1	
9.2	The Godfather	1155770	2	
9	The Dark Knight	1676169	3	
9	The Godfather: Part II	790926	4	
8.9	Pulp Fiction	1324680	5	
8.9	The Lord of the Rings: The Return of the King	1215718	6	
8.9	Schindler's List	865020	7	
8.9	The Good, the Bad and the Ugly	503509	8	

Created top 250 Movies Column in

Insights :

There is top 250 movies according to obtained IMDB Score and sum of user vote I found in this analysis that maximum number (above 20 movies) of movies are in bucket of 7.9 IMDB Score to 8.3.

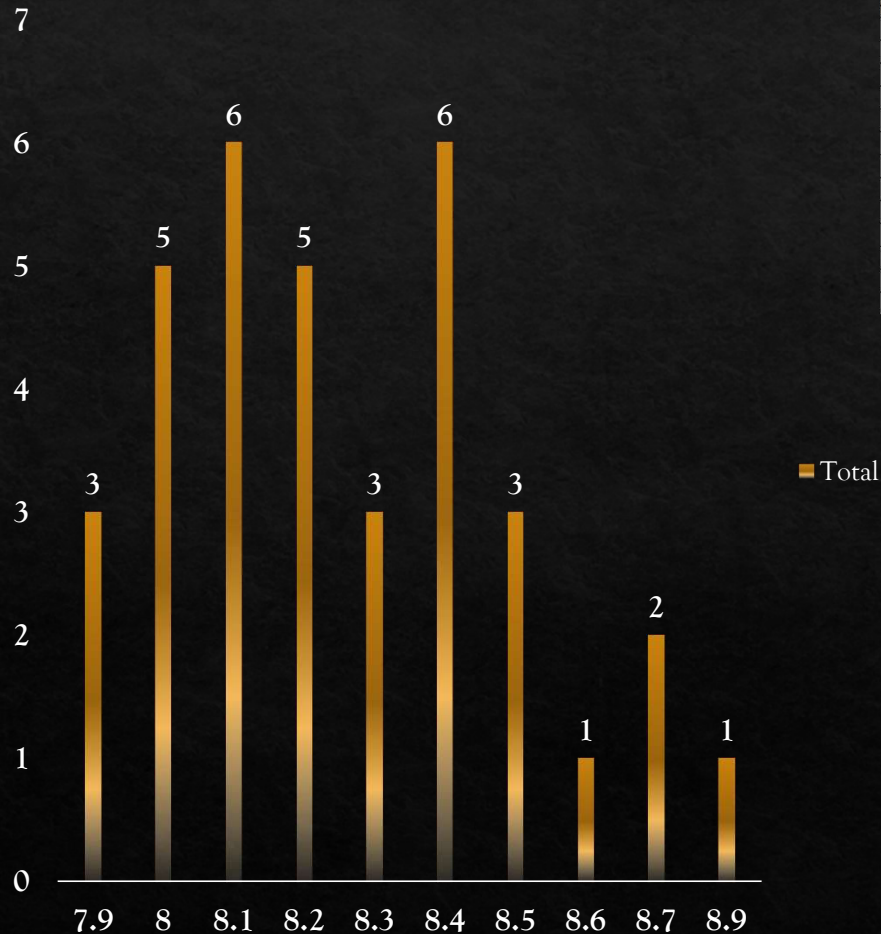


Use these column to create pivot chart and then create bar chart to visualize the data.



# Task C - Top\_Foreign\_Lang\_Film Find into IMDB Top 250

## TOP\_FOREIGN\_LANG\_FILM

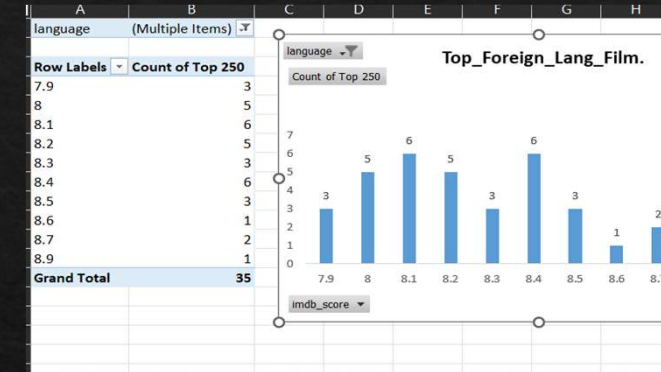


Z	AA	AB	AC	
imdb_score	aspect_ratio	movie_facebook_likes	Profit	Top 250
9.3	1.85	108000	3341469	The Shawshank Redemp
9.2	1.85	43000	128821952	The Godfather
9	2.35	37000	348316061	The Dark Knight
9	1.85	14000	44300000	The Godfather: Part II
8.9	2.35	45000	99930000	Pulp Fiction
8.9	2.35	16000	283019252	The Lord of the Rings: T
8.9	1.85	41000	74067179	Schindler's List
8.9	2.35	20000	4900000	The Good, the Bad and t
8.8	2.35	175000	132568851	Inception
8.8	2.35	48000	-25976605	Fight Club
8.8	2.35	59000	274691196	Forrest Gump
8.8	2.35	21000	220837577	The Lord of the Rings: T
8.8	2.35	17000	272158751	Star Wars: Episode V - T

Created top 250 Movies Column in

Insights :

There is top 250 movies according to obtained IMDB Score and sum of user vote I found in this analysis that the total **35** foreign language movies in this list and maximum number (above 5 movies) of movies belongs to 8.1 & 8.4 IMDB Score.



Use these column to create pivot chart, in the pivot chart I put filter on language column and then create Stacked chart to visualize the data.



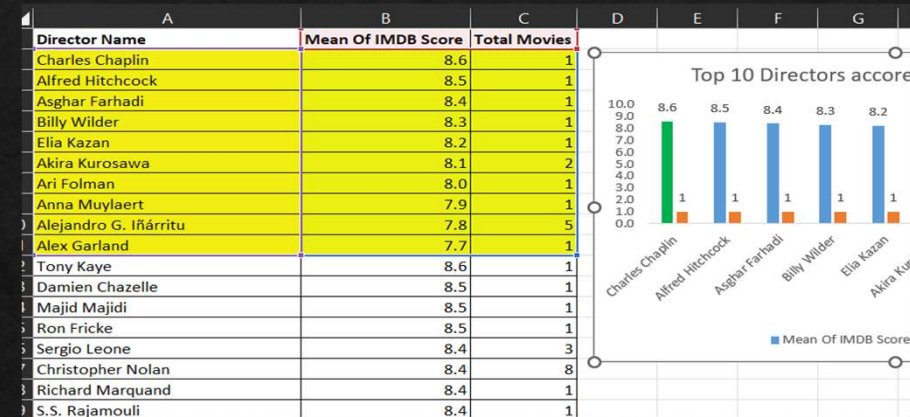
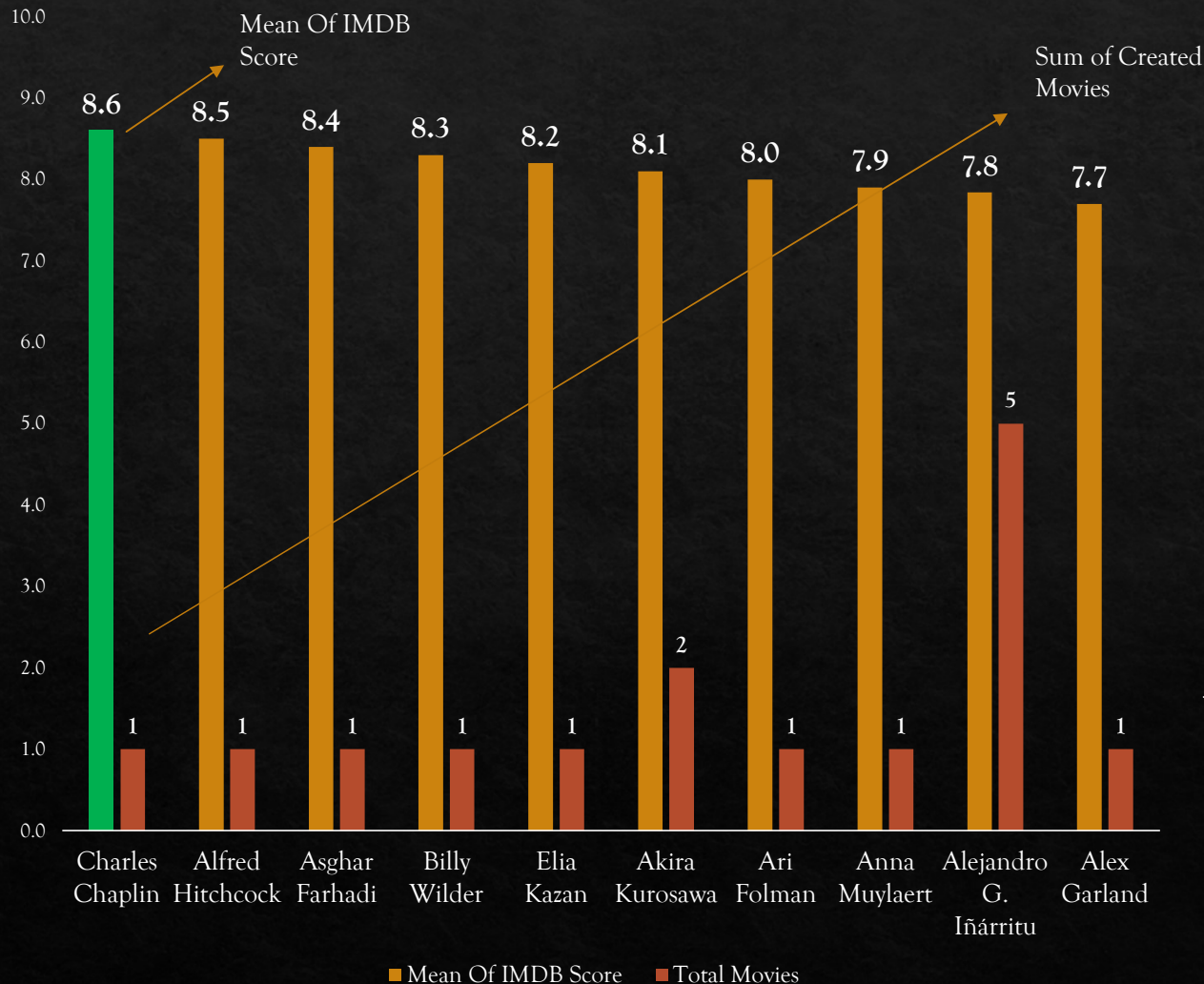
# Task D : Best Directors

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

# Task D - TOP 10 DIRECTORS

Top 10 Directors according to IMDB Score



Apply pivot chart on given dataset, in the pivot chart I put directors name on rows and IMDB score on values as mean of IMDB Score and also put total movies on values the short IMDB Score on the basis of largest to smallest and then I did some conditional formatting for figure out top 10 directors and then create a clustered chart to visualize the data.

## Insights :

There is Top 10 Directors name along with the mean of IMDB Score and Total Created Movies Mr. Charles Chaplin obtained the maximum IMDB Score

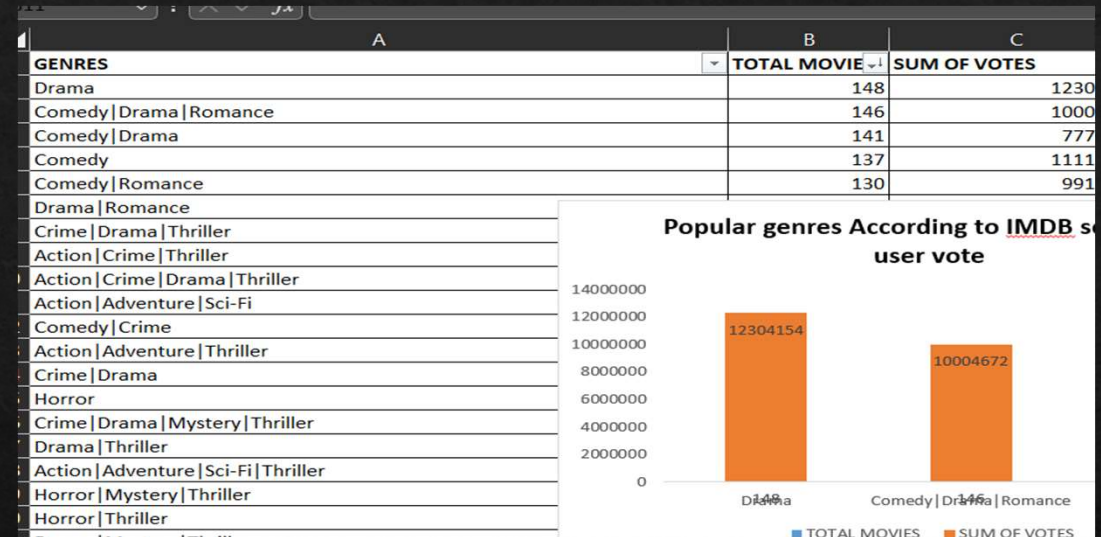
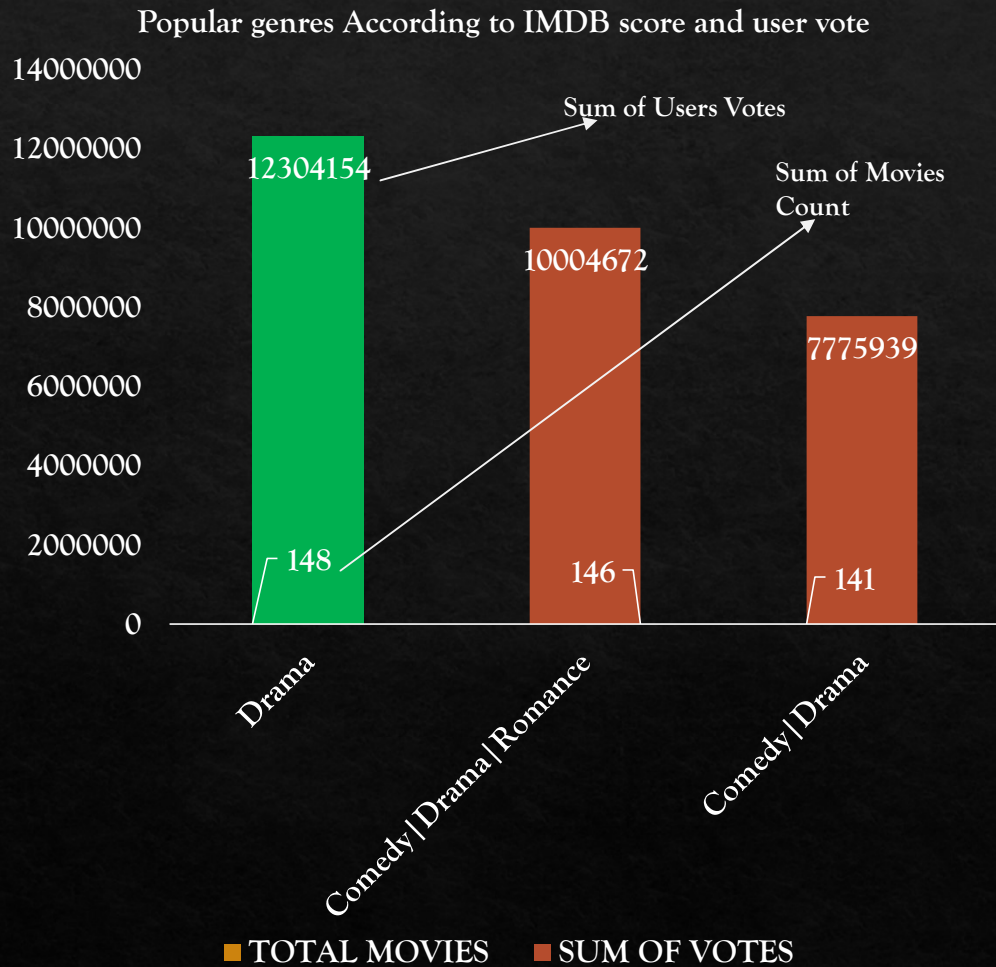
# Task E : Popular Genres

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

Your task : Find popular genres



# Task E - Find the Popular Genres



Apply pivot chart on given dataset, in the pivot chart I put filter large to small on genres count as well as apply same on num\_voted\_users column and then create Stacked chart to visualize the data.

## Insights :

There is Top 3 Genres as it shown in chart the most popular Genres is Drama which is use 148 times as well as the most voted Genres by users is also Drama “The Sum of Users Vote is 12304154 which is the highest one.”

# Task F : Charts

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.

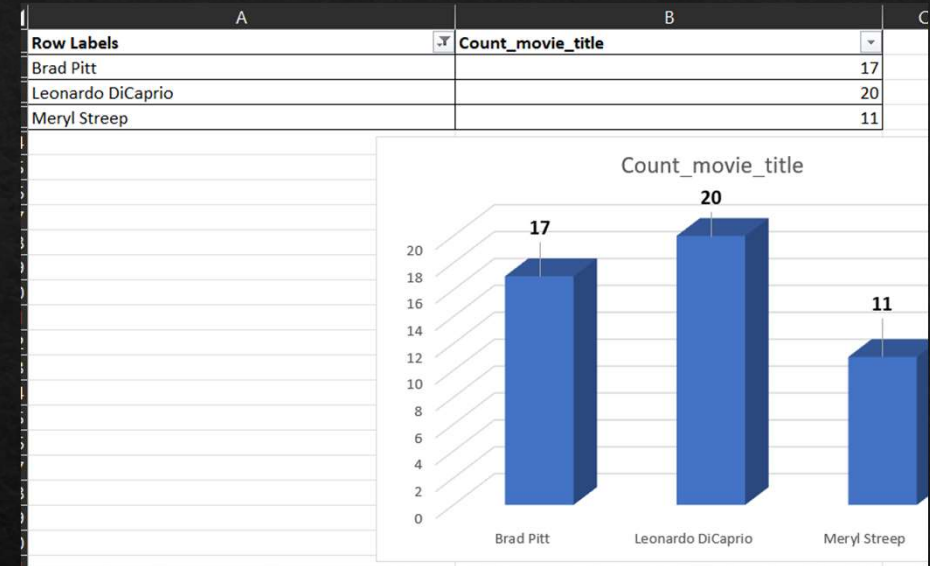
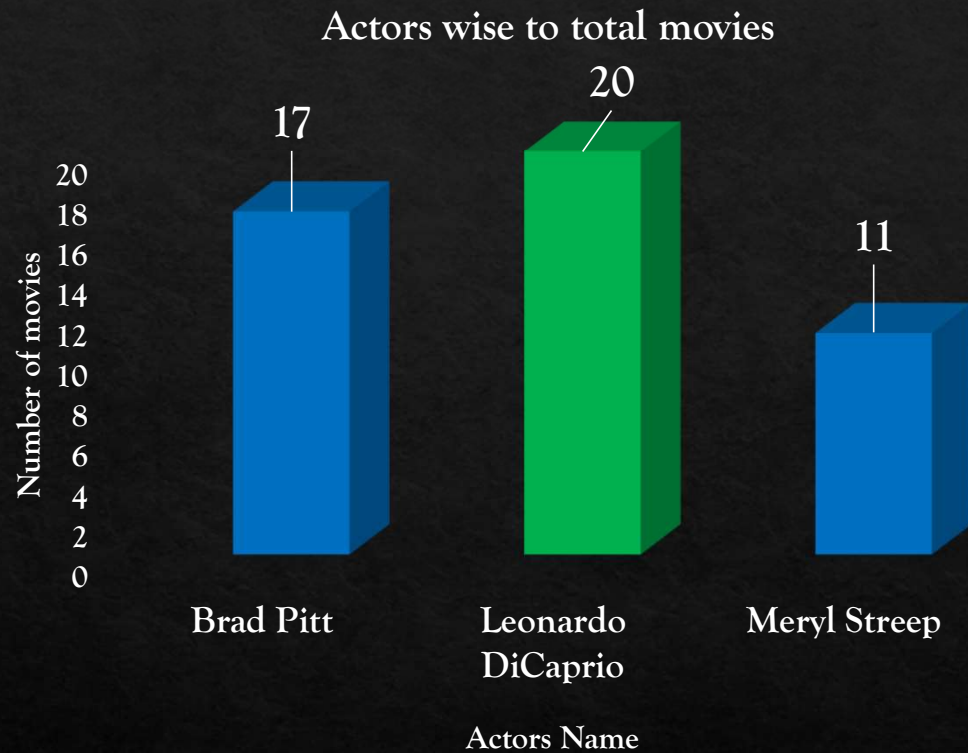
Group the combined column using the actor\_1\_name column.

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

# Task F - Actors wise to total movies



Apply pivot chart on given dataset, in the pivot chart I put filter large to small on count\_movie\_title as well as apply filter by given actors name on actor\_1\_name column and then create 3D Stacked chart to visualize the data.

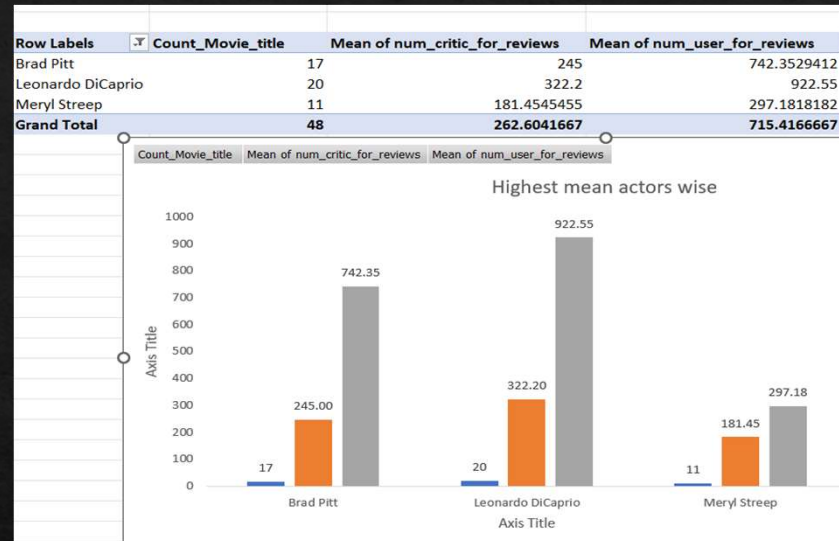
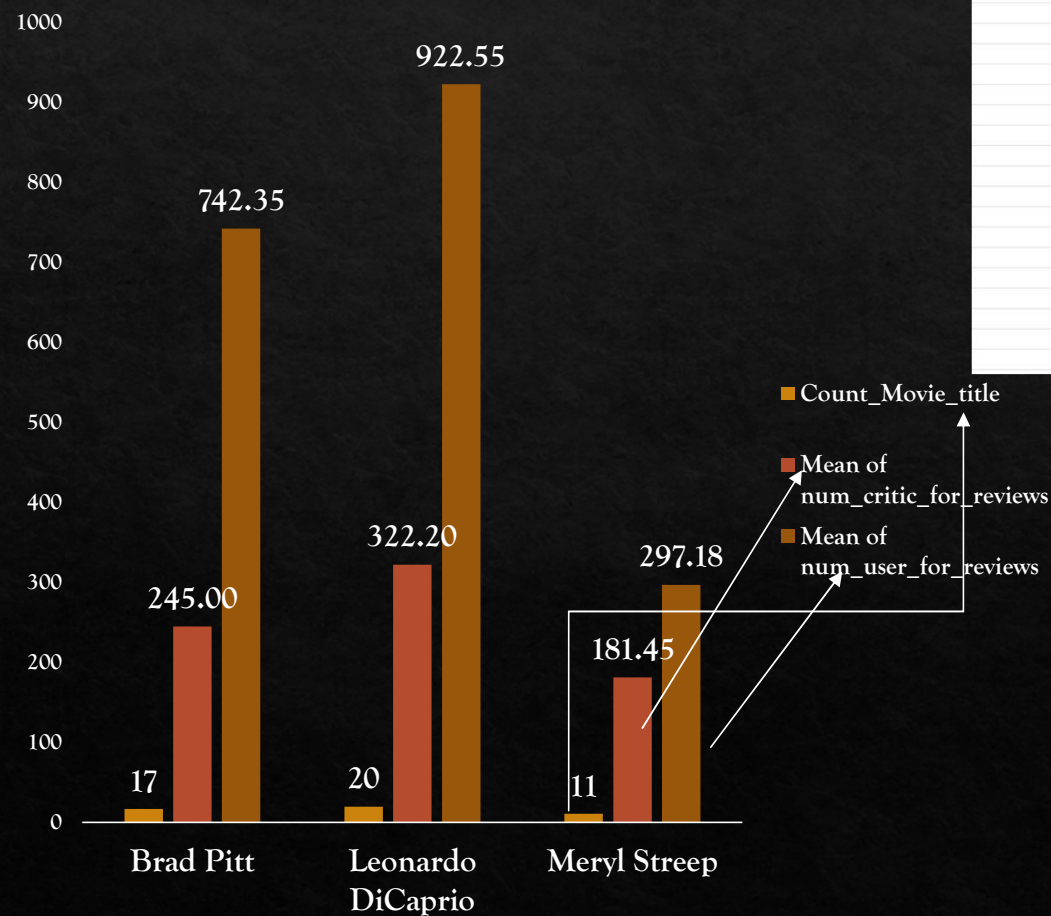
## Insights :

There is 03 Actors with his total movies as per the analysis the chart I found Leonardo DiCaprio worked as a lead actor in 20 movies as well as Bred Pitt 17 and Meryl Streep 11 Hence the max movies done tag goes to Leonardo DiCaprio.



# Task F - Actors wise highest mean

Actors wise mean of  
the num\_critic\_for\_reviews and num\_users\_  
for\_review

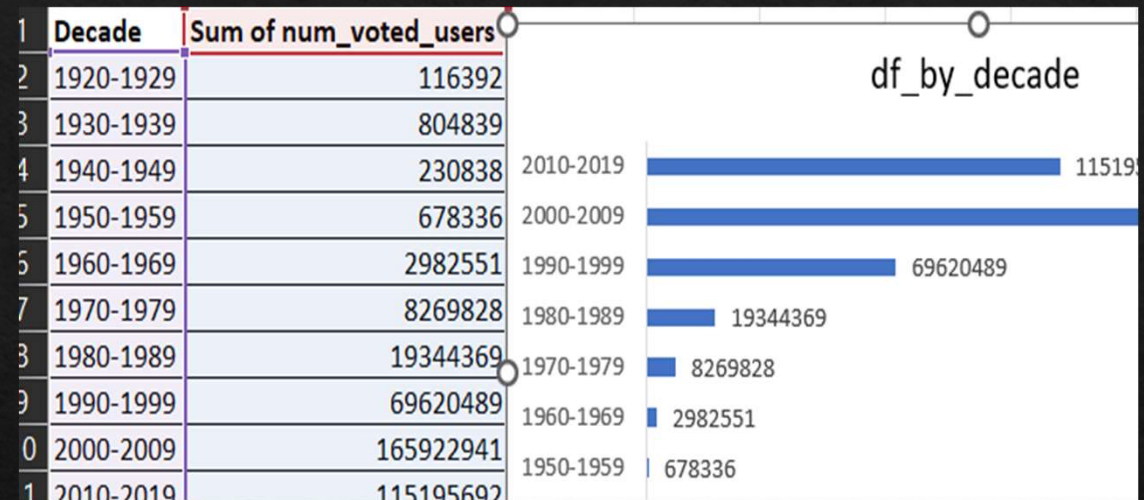
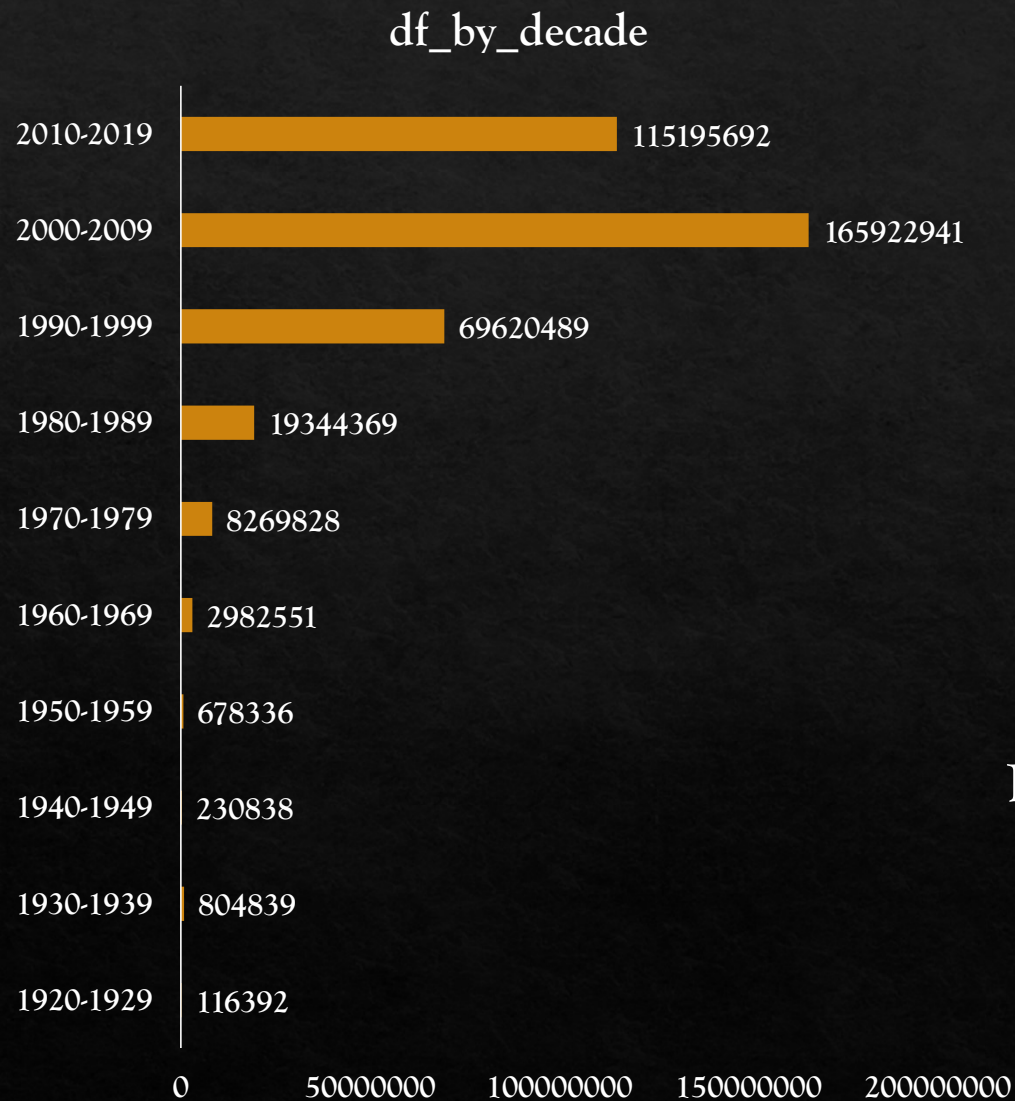


Apply pivot chart on  
extracted data, in the  
pivot chart I take  
mean of critic review  
and mean of user  
review in respective  
of actor name and  
then create clustered  
chart to visualize the  
data.

## Insights :

There is 03 Actors with his total movies, the Mean of critic review and mean of user review as per the analysis the chart I found Leonardo DiCaprio worked as a lead actor in 20 movies and also he get roundoff 922 users review and roundoff 322 critic review as well as Bred Pitt 17 movies and also he get roundoff 742 users review and roundoff 245 critic review and Meryl Streep 11 movies and also Meryl Streep get roundoff 297 users review and roundoff 181 critic review Hence the most favourite actor tag goes to Leonardo DiCaprio.

# Task F - number of voted users over decades using a bar chart



Apply pivot chart on given dataset, in the pivot chart I put title years in rows and Sum of num\_voted\_users in values and group years by 10 Ex- 1920-1929 then I found 10 decades and also it's voted user value to use this value I created a Bar chart to visualize the data.

## Insights :

There is 10 Decades in given bar chart to analysis the chart I found the maximum number of users are voted between year 2000 to 2009 in this decades number of votes is 165922941 and also the minimum number of votes obtained year between 1920 to 1929 the figure of votes is 116392



Thank You !