



IMDB Movie Analysis By Darshan Hedgire

Trainity

IMDB Movie Analysis

Description:

We are provided with dataset having various columns of different IMDB Movies. We are required to Frame the problem. For this task, we will need to define a problem you want to shed some light on.

We can do this by asking 'What?' This is where you frame the problem i.e. What is the problem?

We can do this by asking the following what questions

- What do you see happening?
- What is your hypothesis for the cause of the problem? (this will be broadly based on intuition initially)
- What is the impact of the problem on stakeholders?
- What is the impact of the problem not being solved?

How we are going to handle the things

- Cleaning the data
- Data Analysis skills to explore the data set
- Deriving insights

The things that we are going to find out through the project are movies with the highest profit, top movies as per imdb rating, top directors, most popular genres, top foreign language films and more.

Approach:

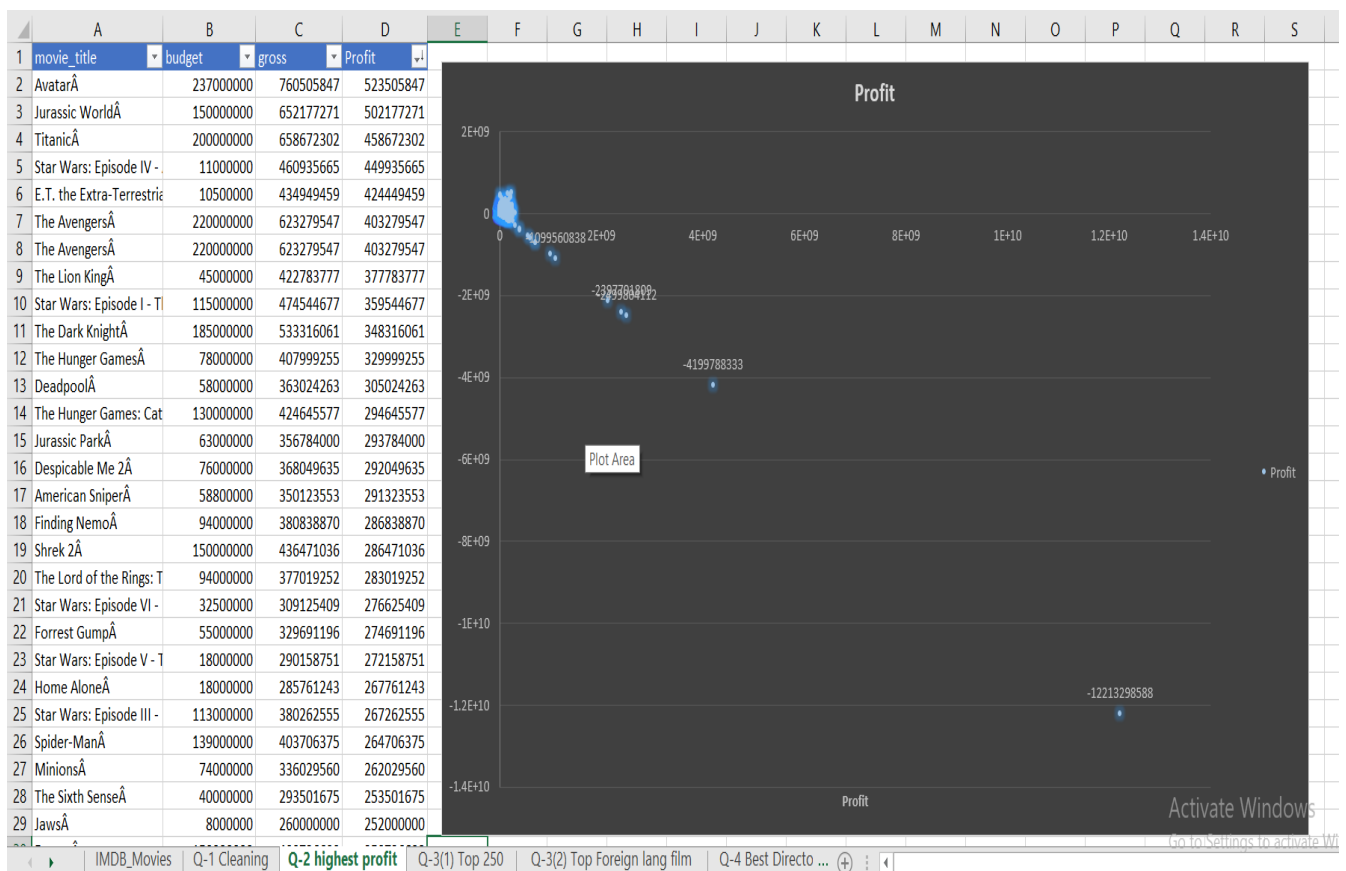
1. Task: Clean the data

This is one of the most important step to perform before moving forward with the analysis.

- First, we dropped the columns which have no use for the analysis.
- Second, we dropped the rows which are blank/null.
- Third, removed the duplicate row values.

2. Task: Find the movies with the highest profit?

- First, we created a new column 'Profit' by subtracting budget column from gross column.
- Second, we sorted the columns using the profit column as reference from the largest to the smallest.
- Third, we plotted the budget and profit in XY Scatter chart to find the outliers.
- There are as many as 5 outliers in the profit columns.
- The movie with the highest profit is 'Avatar' followed by 'Jurassic World' and 'Titanic' and so on.



3. Task: Find IMDB Top 250

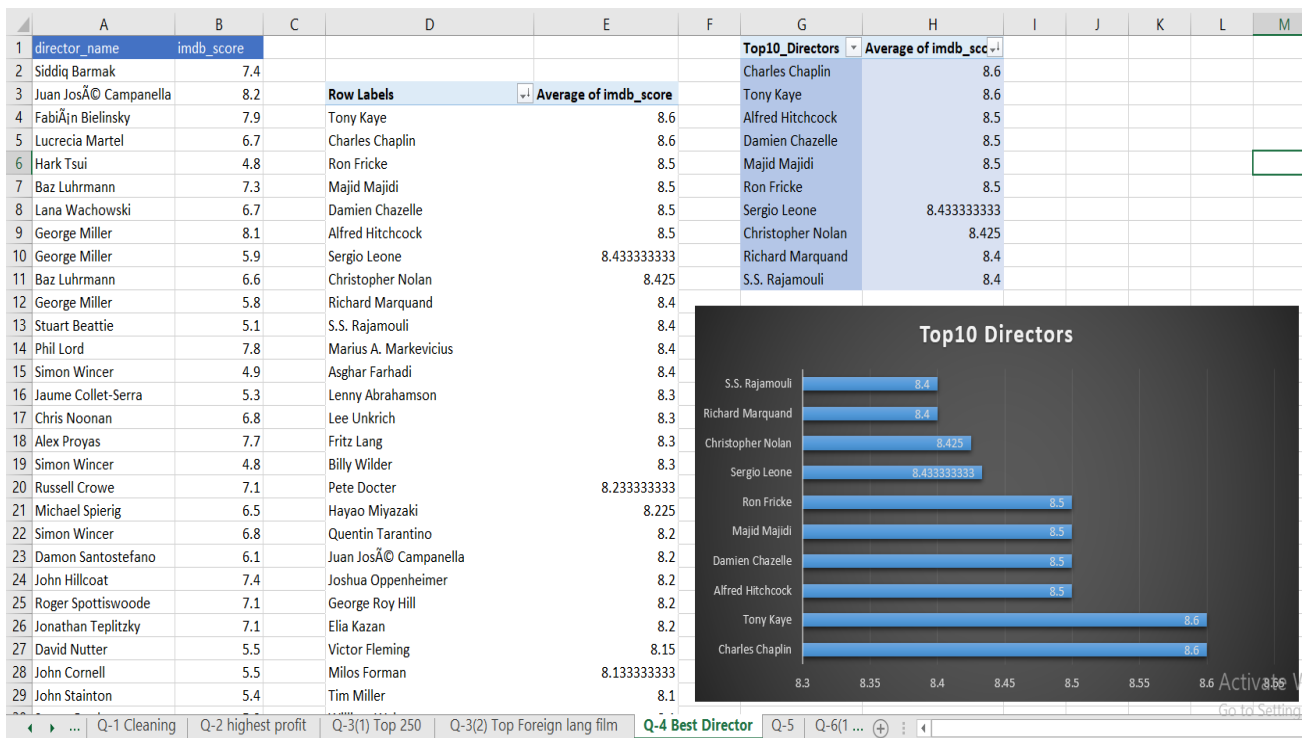
- First, we filtered the 'num_voted_users' column greater than 25,000.
- Second, we created a new column named 'IMDb_Top_250' and stored the top 250 movies with the highest IMDb Rating (sorted the 'imdb_score' column from the largest to the smallest).
- Third, added a 'Rank' containing the values 1 to 250 using the RANK() function + COUNTIFS() function.
- Fourth, we extracted all the movies in the IMDb_Top_250 column by filtering the 'language' column (unselecting English language) and stored them in a new column named 'Top_Foreign_Lang_Film'.

	A	B	C	D	E
1	IMdb_Top_250	num_voted_users	language	imdb_score	Rank
2	The Shawshank Redemption	1689764	English	9.3	1
3	The Godfather	1155770	English	9.2	2
4	The Dark Knight	1676169	English	9	3
5	The Godfather: Part II	790926	English	9	4
6	Pulp Fiction	1324680	English	8.9	5
7	The Lord of the Rings: The Return of the King	1215718	English	8.9	6
8	Schindler's List	865020	English	8.9	7
9	The Good, the Bad and the Ugly	503509	Italian	8.9	8
10	Inception	1468200	English	8.8	9
11	Fight Club	1347461	English	8.8	10
12	Forrest Gump	1251222	English	8.8	11
13	The Lord of the Rings: The Fellowship Ring	1238746	English	8.8	12
14	Star Wars: Episode V - The Empire Strikes Back	837759	English	8.8	13
15	The Matrix	1217752	English	8.7	14
16	The Lord of the Rings: The Two Towers	1100446	English	8.7	15
17	Star Wars: Episode IV - A New Hope	911097	English	8.7	16
18	Goodfellas	728685	English	8.7	17
19	One Flew Over the Cuckoo's Nest	680041	English	8.7	18
20	City of God	533200	Portuguese	8.7	19
21	Seven Samurai	229012	Japanese	8.7	20
22	Se7en	1023511	English	8.6	21
23	Interstellar	928227	English	8.6	22
24	The Silence of the Lambs	887467	English	8.6	23
25	Saving Private Ryan	881236	English	8.6	24
26	American History X	782437	English	8.6	25
27	The Usual Suspects	740918	English	8.6	26
28	Spirited Away	417971	Japanese	8.6	27
29	Modern Times	143086	English	8.6	28

	A	B	C	D
1	Top_Foreign_Lang_Film	num_voted_users	language	imdb_score
2	The Good, the Bad and the Ugly	503509	Italian	8.9
3	City of God	533200	Portuguese	8.7
4	Seven Samurai	229012	Japanese	8.7
5	Spirited Away	417971	Japanese	8.6
6	The Lives of Others	259379	German	8.5
7	Children of Heaven	27882	Persian	8.5
8	Amélie	534262	French	8.4
9	Oldboy	356181	Korean	8.4
10	Princess Mononoke	221552	Japanese	8.4
11	Das Boot	168203	German	8.4
12	A Separation	151812	Persian	8.4
13	Baahubali: The Beginning	62756	Telugu	8.4
14	Downfall	248354	German	8.3
15	The Hunt	170155	Danish	8.3
16	Metropolis	111841	German	8.3
17	Pan's Labyrinth	467234	Spanish	8.2
18	Howl's Moving Castle	214091	Japanese	8.2
19	The Secret in Their Eyes	131831	Spanish	8.2
20	Incendies	80429	French	8.2
21	Amores Perros	173551	Spanish	8.1
22	Akira	106160	Japanese	8.1
23	Elite Squad	81644	Portuguese	8.1
24	The Celebration	65951	Danish	8.1
25	The Sea Inside	64556	Spanish	8.1
26	Tae Guk Gi: The Brotherhood of War	31943	Korean	8.1
27	A Fistful of Dollars	147566	Italian	8
28	Persepolis	70194	French	8
29	My Name Is Khan	69759	Hindi	8
30	My Name Is Khan	69759	Hindi	8

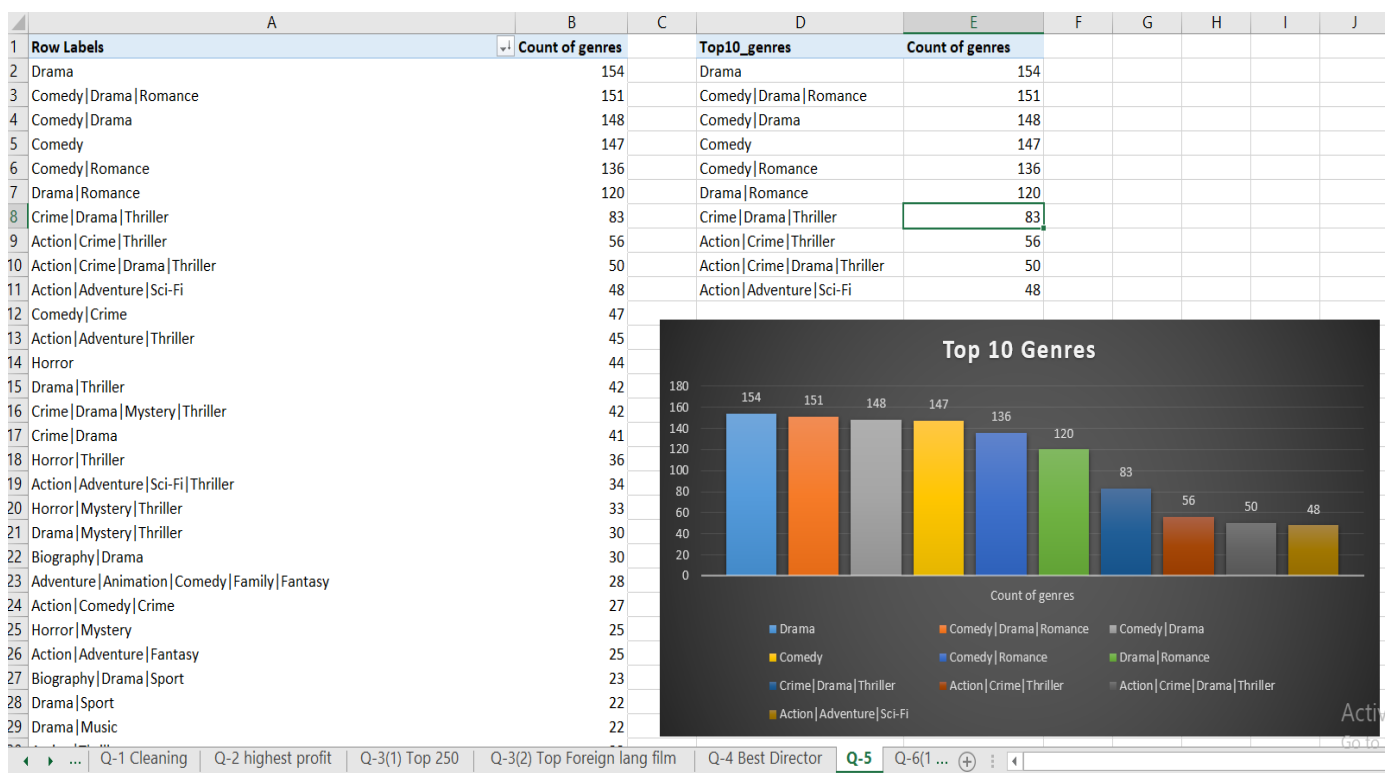
4. Task: Find the best directors

- First, we selected the cleaned dataset done in Task 1 and created a pivot table.
- Second, we put the 'director_name' into the Rows and took average of 'imdb_score' in the Values section.
- Third, we sorted the 'director_name' in ascending order and then sorted the 'average of imdb_score' (largest to smallest).
- Then we selected the top 10 directors and their mean of imdb_score in other columns.
- Next, we made a bar chart of the top 10 directors for the better insights.



5. Task: Find popular genres

- First, we selected the 'genres' column from the cleaned dataset done in Task 1 and created a pivot table.
- Second, we put the 'genres' into the Rows and took count of 'genres' in the Values section.
- Third, we sorted the 'Count of genres' in descending order.
- Then we copied the top 10 genres and their count and pasted it in the other columns.
- Next, we made a Column chart of the top 10 genres for the better insights.



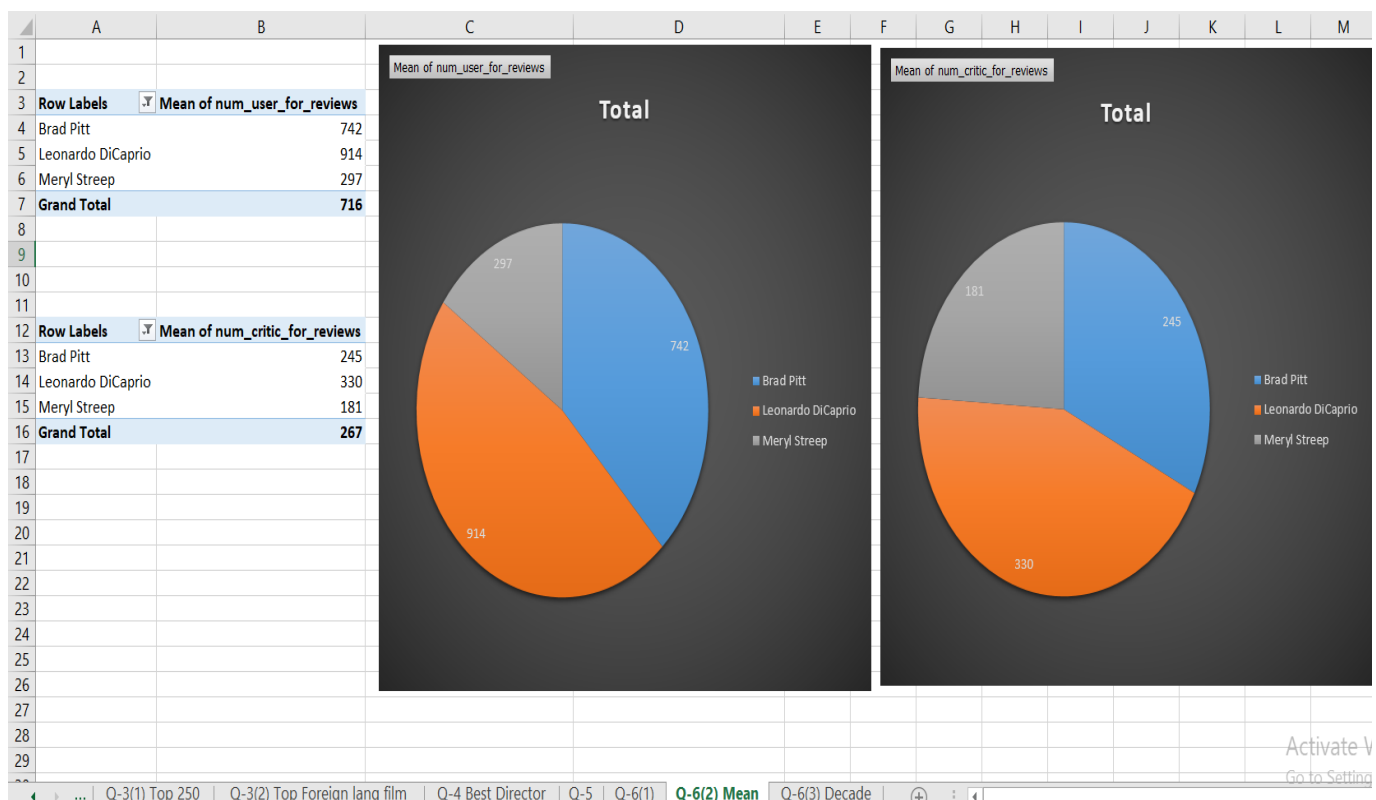
6. (1) Task: Find the critic-favorite and audience-favorite actors

- First, we created 3 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 using the 'actor_1_name' column.
- Second, appended the rows of all these columns and stored them in a new column named 'Combined'.
- We grouped the column by the actor's name: 'Meryl Streep', 'Leonardo DiCaprio', and 'Brad Pitt'.

	C	D	E	F	G	H	I	J	K
1	num_critic_for_review	num_user_for_review	num_voted_users	title_year	Meryl_Streep	Leo_Caprio	Brad_Pitt	Combined	Group_By
2	105	77	7559	2003	The Iron Lady	The Great Gatsby	Babel	The Iron Lady	Meryl Streep
3	262	231	131831	2009	It's Complicated	The Great Gatsby	Spy Game	It's Complicated	Meryl Streep
4	94	125	38215	2000	The River Wild	Blood Diamond	The Curious Case of Benjamin Button	The River Wild	Meryl Streep
5	78	37	2720	2004	Julie & Julia	The Quick and the Dead	Troy	Julie & Julia	Meryl Streep
6	67	141	11512	1998	The Devil Wears Prada	Titanic	Ocean's Twelve	The Devil Wears Prada	Meryl Streep
7	490	753	362912	2013	Lions for Lambs	Inception	Mr. & Mrs. Smith	Lions for Lambs	Meryl Streep
8	245	2121	364948	2003	Out of Africa	The Revenant	Ocean's Eleven	Out of Africa	Meryl Streep
9	739	1588	552503	2015	Hope Springs	The Aviator	Fury	Hope Springs	Meryl Streep
10	156	79	32399	2011	One True Thing	Django Unchained	Seven Years in Tibet	One True Thing	Meryl Streep
11	245	450	102338	2008	The Hours	The Wolf of Wall Street	Fight Club	The Hours	Meryl Streep
12	61	216	24868	1998	A Prairie Home Companion	Gangs of New York	Sinbad: Legend of the Seven Seas	A Prairie Home Companion	Meryl Streep
13	308	194	65709	2014		The Departed	Interview with the Vampire: The Vampire Chronicles	The Great Gatsby	Leonardo DiCaprio
14	435	471	246698	2014		Shutter Island	The Tree of Life	The Great Gatsby	Leonardo DiCaprio
15	60	148	27648	1996		Body of Lies	The Assassination of Jesse James by the Coward Jim Bowie	Blood Diamond	Leonardo DiCaprio
16	228	576	94456	2005		Catch Me If You Can	By the Sea	The Quick and the Dead	Leonardo DiCaprio
17	79	182	94435	1995		The Beach	Killing Them Softly	Titanic	Leonardo DiCaprio
18	222	624	156929	1998		Revolutionary Road	True Romance	Inception	Leonardo DiCaprio
19	71	119	19699	2001		The Man in the Iron Mask		The Revenant	Leonardo DiCaprio
20	183	185	53341	2014		J. Edgar		The Aviator	Leonardo DiCaprio
21	292	281	105797	2009		Marvin's Room		Django Unchained	Leonardo DiCaprio
22	20	97	15230	1990		Romeo + Juliet		The Wolf of Wall Street	Leonardo DiCaprio
23	81	121	16372	1999				Gangs of New York	Leonardo DiCaprio
24	168	232	43205	2005				The Departed	Leonardo DiCaprio
25	67	33	8087	2008				Shutter Island	Leonardo DiCaprio
26	186	119	27882	2013				Body of Lies	Leonardo DiCaprio
27	96	237	17328	1998				Catch Me If You Can	Leonardo DiCaprio
28	32	62	44096	1988				The Beach	Leonardo DiCaprio
29	61	106	5663	2002				Revolutionary Road	Leonardo DiCaprio
30	58	68	23383	1992				The Man in the Iron Mask	Leonardo DiCaprio
31	138	354	18632	2003				J. Edgar	Leonardo DiCaprio
32	25	12	1427	1997				Marvin's Room	Leonardo DiCaprio
33	75	115	5772	2006				Romeo + Juliet	Leonardo DiCaprio
34	117	193	101840	1985				Babel	Brad Pitt
35	24	6	813	2006				Spy Game	Brad Pitt
36	35	93	74743	1986				The Curious Case of Benjamin Button	Brad Pitt
37	21	43	4792	1989				Troy	Brad Pitt

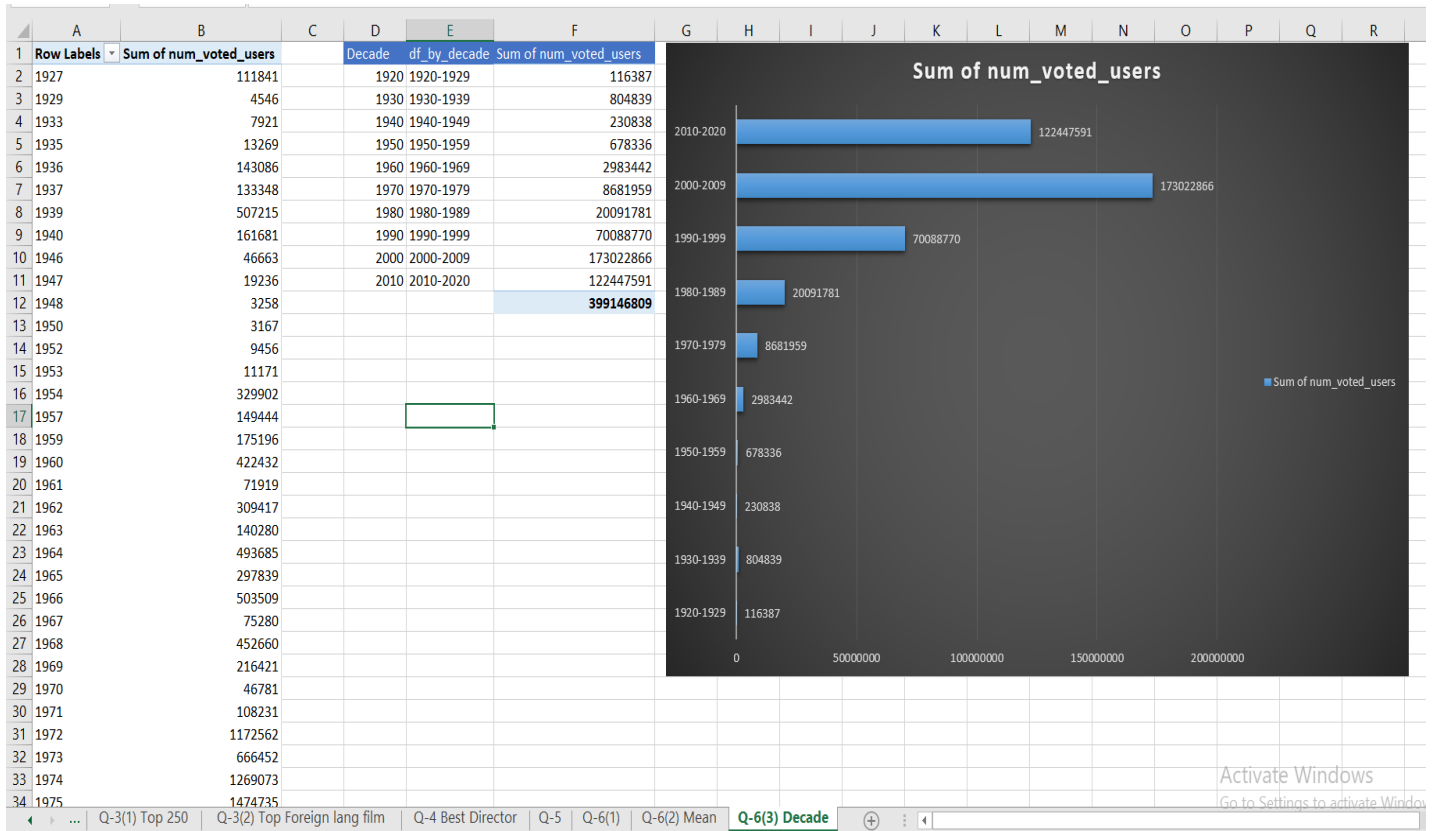
6 (2) Find the mean

- We selected the cleaned dataset done in Task 1 and created a pivot table.
- Next, we put the 'actor_1_name' into the Rows and took mean/average of 'num_users_for_review' in the Values section.
- Sorted the column from largest to smallest by mean of 'num_users_for_review'.
- Then, we made a pivot chart (pie chart) of the mean of 'num_users_for_review'.
- Again, we did the same above process for the mean of 'num_critic_for_review'.



6 (3) Change in number of voted users over decades:

- First, we selected the cleaned dataset done in Task 1 and created a pivot table.
- Second, we put the 'title_year' into the Rows and took the sum of 'num_voted_users' in the Values section.
- Third, we grouped the title_year by decade and stored in df_by_decade column.
- Lastly, we plotted the total no. of voted users against the decade in a bar chart.



Tech-Stack Used:

- Microsoft Excel 2016: It enables users to format, organize and calculate data in a spreadsheet. It organizes data in an easy-to-navigate way. We need not to perform any complex mathematical functions. And it turns piles of data into helpful graphics and charts.
- Microsoft Word 2016: It is used to make a report (PDF) to be presented to the leadership team.

Insights:

- There are as many as 5 outliers in the profit columns.
- The movie with the highest profit is 'Avatar' followed by 'Jurassic World' and 'Titanic' and so on.
- The Shawshank Redemption is the top-most movie with the highest IMDB rating.
- The Good, the Bad and the Ugly (Italian) is the top-most foreign language movie.
- Charles Chaplin is the top-most director followed by Tony Kaye.
- The most popular genres is Drama followed by Comedy.
- 'Leonardo DiCaprio' is the critic-favorite as well as the audience-favorite actor.
- The most users voted in the decade 2000s and the least in the decade 1940s.

Results:

- In this project, I applied the basic and advance Excel concepts. The concepts related to statistics and EDA have been implemented here by using MS Excel.
- In this task, the concepts regarding the sort, filter, pivot table, charts, different functions like rank, etc have been implemented.
- I learned to implement the learning of Excel in the real-time project.
- I learned how to frame the problem by asking 'what' looking at the dataset.
- It helped me in learning the '5 Why Analysis' to determine the root cause of the problem.
- I learned how a data analyst think deeper and deeper to generate the valuable insights.
- It was a great learning experience while doing this project and it was challenging too while asking the different questions and finding their answers.

Excel Sheet Link:

<https://drive.google.com/drive/folders/1YDvsRlae9L1LYISLdN5hJ7kAnxFzDN7m?usp=sharing>