

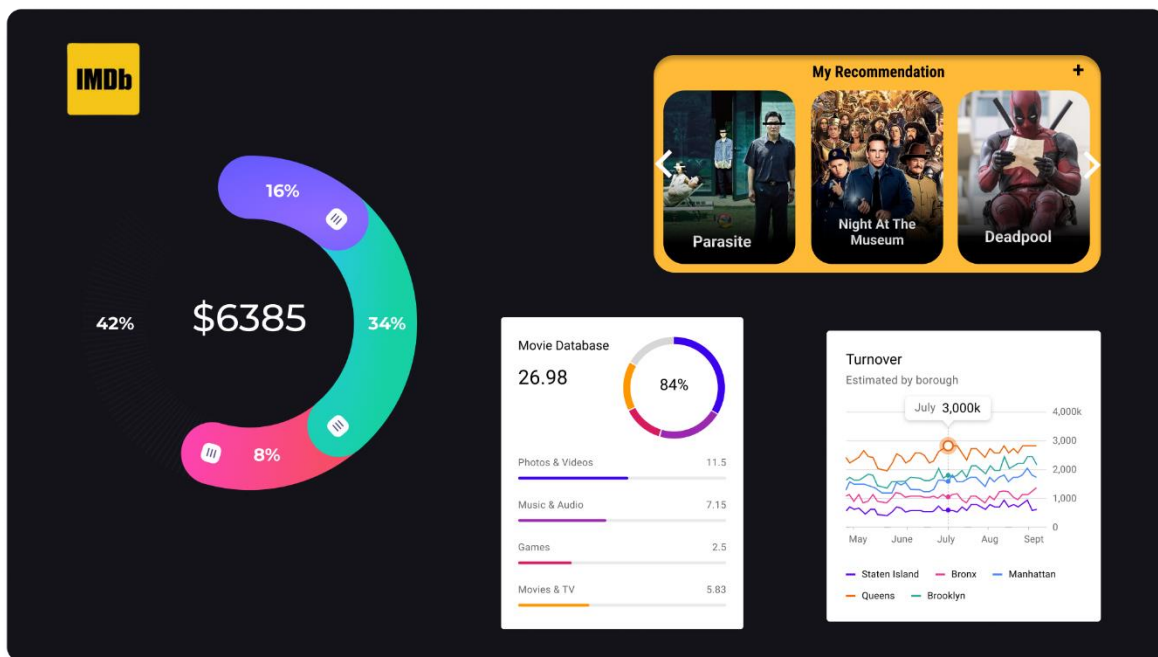
# IMDB Movie Analysis

(Final Project-1)

By Shahequa Modabbera

[LinkedIn](#)

[GitHub](#)



## Project Description:

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

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

We can do this by asking the following 'What?' :

- What do we see happening?
- What is our 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 to handle the things:

- Clean the data.

- Use the Data Analysis skills to explore the data set.
- Derive 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:

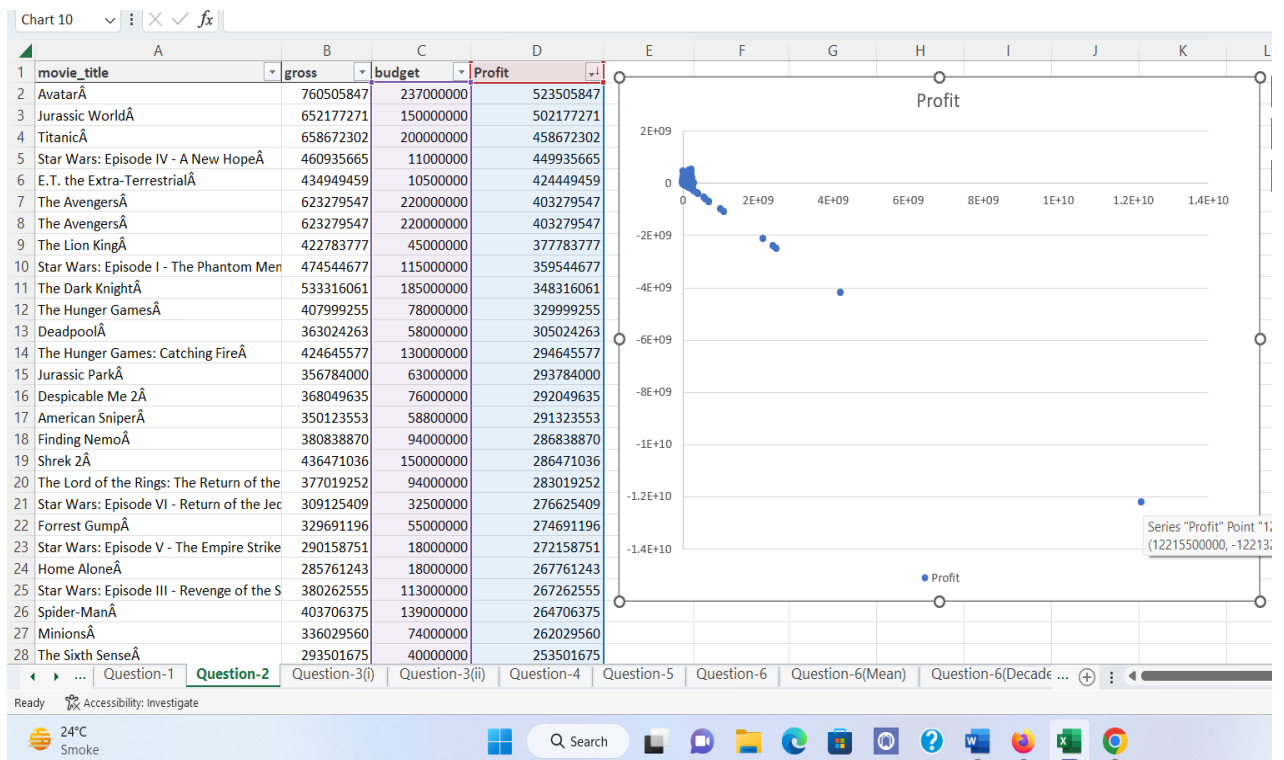
### 1. Task: Clean the data

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, 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, 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	Schindler's List	865020	English	8.9	6
8	The Good, the Bad and the Ugly	503509	Italian	8.9	7
9	The Lord of the Rings: The Return of the King	1215718	English	8.9	8
10	Fight Club	1347461	English	8.8	9
11	Forrest Gump	1251222	English	8.8	10
12	Inception	1468200	English	8.8	11
13	Star Wars: Episode V - The Empire Strikes Back	837759	English	8.8	12
14	The Lord of the Rings: The Fellowship of the Ring	1238746	English	8.8	13
15	City of God	533200	Portuguese	8.7	14
16	Goodfellas	728685	English	8.7	15
17	One Flew Over the Cuckoo's Nest	680041	English	8.7	16
18	Seven Samurai	229012	Japanese	8.7	17
19	Star Wars: Episode IV - A New Hope	911097	English	8.7	18
20	The Lord of the Rings: The Two Towers	1100446	English	8.7	19
21	The Matrix	1217752	English	8.7	20
22	American History X	782437	English	8.6	21
23	Interstellar	928227	English	8.6	22
24	Modern Times	143086	English	8.6	23
25	Saving Private Ryan	881236	English	8.6	24
26	Seven Years in Tibet	1023511	English	8.6	25
27	Spirited Away	417971	Japanese	8.6	26
28	The Silence of the Lambs	887467	English	8.6	27

◀ ▶ ...
Question-1
Question-2
Question-3(i)
Question-3(ii)
Question-4
Question-5
Q

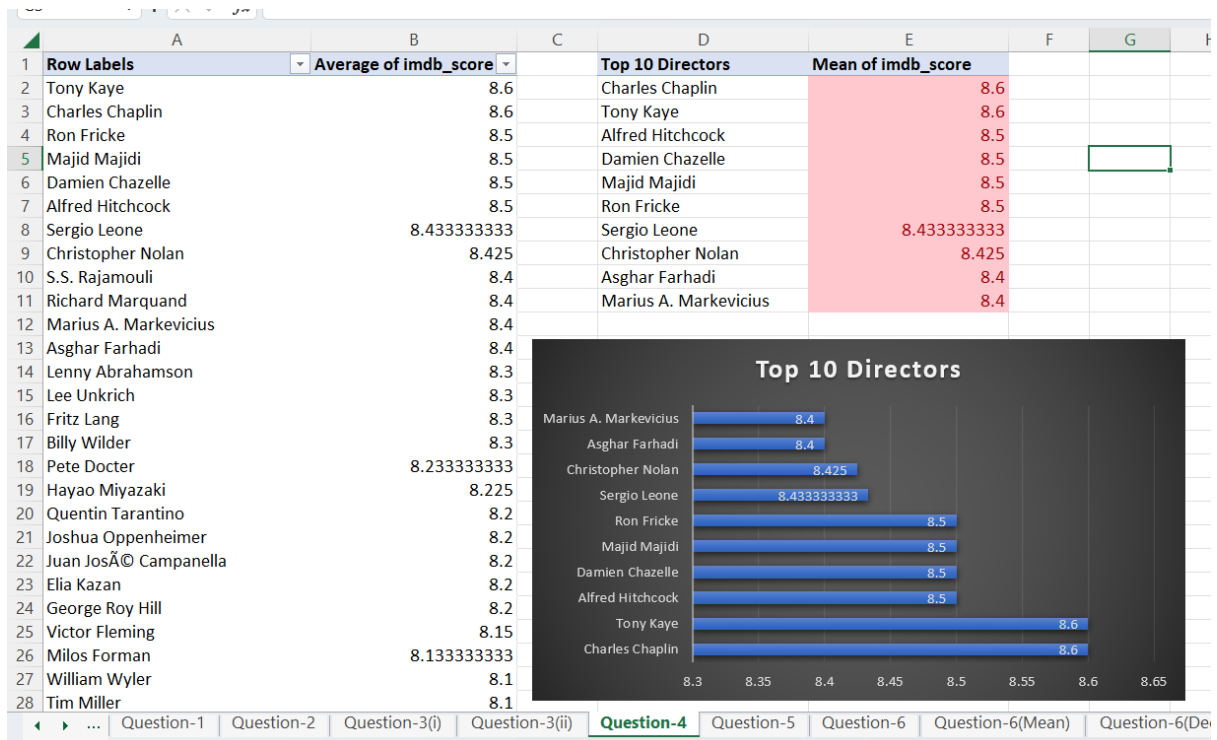
Ready
 Accessibility: Investigate

	A	B	C	D
1	Top_Foreign_Lang_Film	num_voted_users	imdb_score	language
2	The Good, the Bad and the Ugly	503509	8.9	Italian
3	City of God	533200	8.7	Portuguese
4	Seven Samurai	229012	8.7	Japanese
5	Spirited Away	417971	8.6	Japanese
6	Children of Heaven	27882	8.5	Persian
7	The Lives of Others	259379	8.5	German
8	A Separation	151812	8.4	Persian
9	Amélie	534262	8.4	French
10	Baahubali: The Beginning	62756	8.4	Telugu
11	Das Boot	168203	8.4	German
12	Oldboy	356181	8.4	Korean
13	Princess Mononoke	221552	8.4	Japanese
14	Metropolis	111841	8.3	German
15	The Hunt	170155	8.3	Danish
16	Unforgiven	248354	8.3	German
17	Pan's Labyrinth	80429	8.2	French
18	The Bridge on the River Kwai	131831	8.2	Spanish
19	The Thing	467234	8.2	Spanish
20	Warrior	214091	8.2	Japanese
21	Annie Hall	81644	8.1	Portuguese
22	In the Shadow of the Moon	65951	8.1	Danish
23	Sling Blade	106160	8.1	Japanese
24	Tae Guk Gi: The Brotherhood of War	64556	8.1	Spanish
25	The Best Years of Our Lives	173551	8.1	Spanish
26	The Imitation Game	31943	8.1	Korean
27	Bowling for Columbine	28951	8	Portuguese
28	Jaws	70194	8	French

◀ ▶ ... Question-1 Question-2 Question-3(i) **Question-3(ii)** Question-4

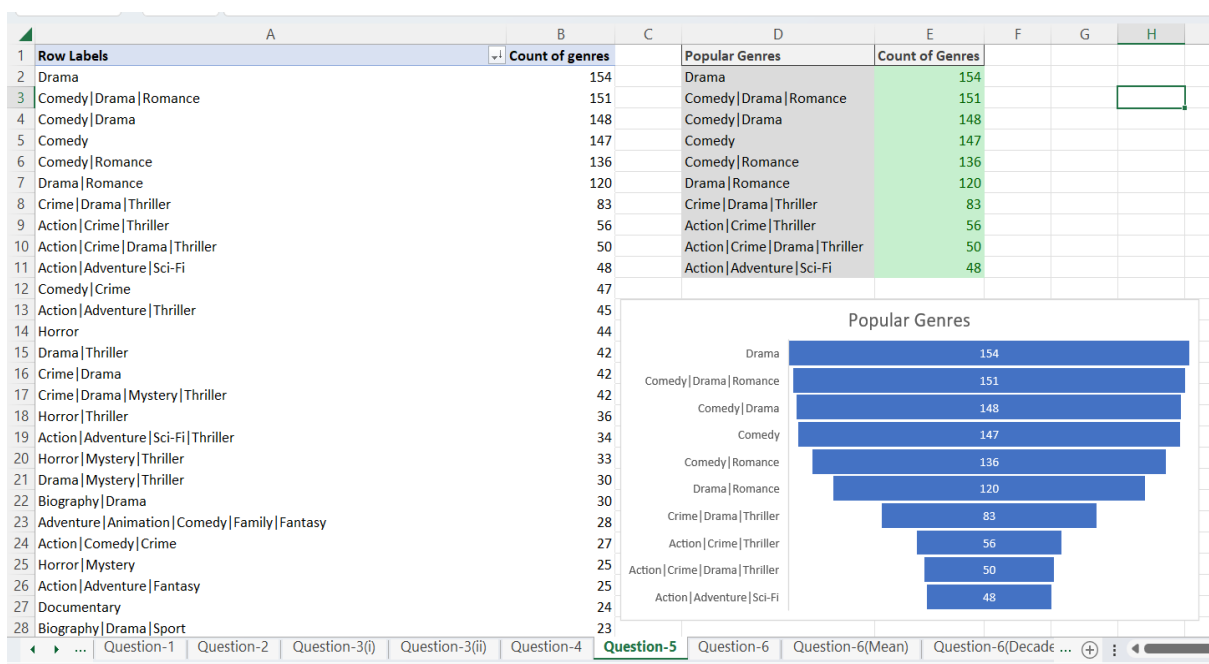
#### 4. Task: Find the best directors

- 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

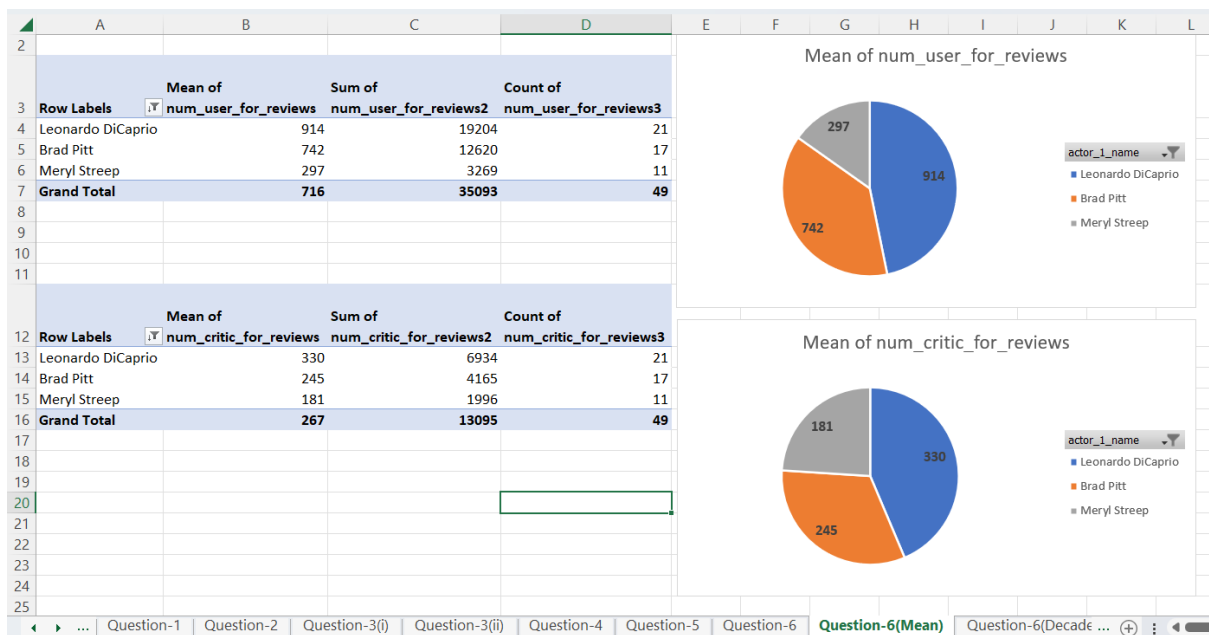
- First, we selected the 'genres' column from the cleaned dataset done in Task 1 and created a pivot table.
- Second, we put the 'genres' into the Rows and took count of 'genres' in the Values section.
- Third, we sorted the 'Count of genres' in descending order.
- Then we copied the top 10 genres and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



## 6. Task: Find the critic-favorite and audience-favorite actors

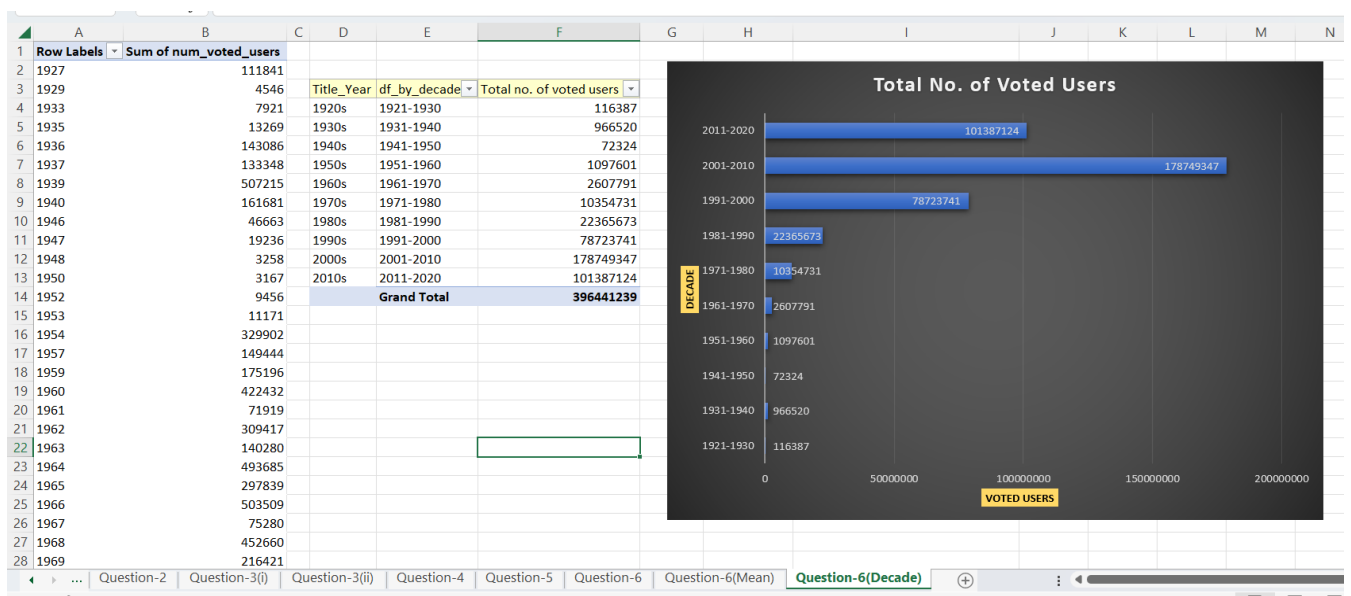
- First, we created 3 new columns namely, Meryl\_Streep, Leo\_Caprio, and Brad\_Pitt which contain the movies in which the actors: 'Meryl Streep', 'Leonardo DiCaprio', and 'Brad Pitt' are the lead actors using the 'actor\_1\_name' column.
- Second, appended the rows of all these columns and stored them in a new column named 'Combined'.
- We grouped the column by the actor's name: 'Meryl Streep', 'Leonardo DiCaprio', and 'Brad Pitt'.
- Then, 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 (bar chart) of the mean of 'num\_users\_for\_review'.
- Again, we did the same above process for the mean of 'num\_critic\_for\_review'.

	B	C	D	E	F	G	H	I	J	K
	num_critic_for_reviews	num_user_for_reviews	num_voted_users	movie_title	title_year	Meryl_Streep	Leo_Caprio	Brad_Pitt	Combined	Group By
1	813	2701	1144337	The Dark Knight Rises	2012	Lucy	Django Unchained	Divergent	Lucy	Meryl Str
2	775	2326	456260	Soul Food	2012	Maleficent	Inception	Shanghai Knights	Maleficent	Meryl Str
3	765	1193	955174	Django Unchained	2012	Transcendence	How to Train Your Dragon	The Circle	Transcendence	Meryl Str
4	750	1498	522048	This Is England	2012	Down and Out with the Dolls	Pirates of the Caribbean: The Curse of the Black Pearl	42nd Street	Down and Out with the Dolls	Meryl Str
5	750	1498	522030	Saw	2012	Machine Gun Preacher	The Sessions	Fight Club	Machine Gun Preacher	Meryl Str
6	739	1588	552503	Guardians of the Galaxy	2015	Like Crazy	Frozen River	Scott Pilgrim vs. the World	Like Crazy	Meryl Str
7	738	1885	582917	Weekend	2013	Brown Sugar	The Martian	Mulan	Brown Sugar	Meryl Str
8	733	2536	548573	Green Room	2013	Bridge of Spies	Lethal Weapon 4	Leap Year	Bridge of Spies	Meryl Str
9	723	3054	886204	Maurice	2009	True Lies	The Departed	The Kids Are All Right	True Lies	Meryl Str
10	712	2725	928227	Interstellar	2014	Hellboy II: The Golden Army	The Hunger Games	Connie and Carla	Hellboy II: The Golden Army	Meryl Str
11	703	1722	995415	A Fistful of Dollars	2012	Gerry	Animal House	Star Trek Into Darkness	Gerry	Meryl Str
12	703	1722	995415	Akira	2012		The Squid and the Whale	Sunshine State	Django Unchained	Leonardo
13	682	678	245333	A Room with a View	2011		Dream with the Fishes	Zookeeper	Inception	Leonardo
14	676	1264	431578	Black Book	2011		Transamerica	Black Swan	How to Train Your Dragon	Leonardo
15	673	3018	371639	Sliding Doors	2016		The Iron Giant	Kung Fu Panda	Pirates of the Caribbean: The Curse of the Black Pearl	Leonardo
16	673	1959	701607	Tarnation	2012		Casino Royale	The Puffy Chair	The Sessions	Leonardo
17	669	1140	551363	In Bruges	2010		Hurricane Streets	It Follows	Frozen River	Leonardo
18	663	269	70336	The Queen	2012		The Devil Wears Prada		The Martian	Leonardo
19	656	695	452465	Cast Away	2012		Ocean's Twelve		Lethal Weapon 4	Leonardo
20	654	995	465019	Sunshine Cleaning	2013		The Longest Yard		The Departed	Leonardo
21	653	1097	682155	The Avengers	2014				The Hunger Games	Leonardo
22	645	1367	637246	Gangster's Paradise: Jerusalem	2012		The Family Man		Animal House	Leonardo
23	645	4667	1676169	The Dark Knight	2008				The Squid and the Whale	Leonardo
24	644	1290	418214	Wild Target	2015				Dream with the Fishes	Leonardo
25	642	2803	1468200	Inception	2010				Transamerica	Leonardo
26	635	1117	462669	Layer Cake	2015				The Iron Giant	Leonardo
27	634	986	277172	Rio	2012				Casino Royale	Leonardo
28	608	1187	557489	Secretary	2013				Hurricane Streets	Leonardo
29	606	1138	780588	How to Train Your Dragon	2013				The Devil Wears Prada	Leonardo
30	602	994	275868	Machete	2015				Ocean's Twelve	Leonardo
31	599	1225	451803	The Rose	2012				The Longest Yard	Leonardo
32	597	695	439176	There Will Be Blood	2012				The Family Man	Leonardo
33	596	1018	272839	Iron Man	2015				Divergent	Brad Pitt
34	590	1171	395573	You Can Count on Me	2013				Shanghai Knights	Brad Pitt
35	590	667	438016	I Found the Virgin of England	2013				The Circle	Brad Pitt



**Second part**(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 365:** It enables users to format, organize and calculate data in a spreadsheet. It organize data in an easy-to-navigate way. We need not to perform any complex mathematical functions. And it turn piles of data into helpful graphics and charts.

- **Microsoft Word 2021:** 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:**

- [IMDB Movie Analysis.xlsx](#)
- Make sure to open it in the MS Excel for the right visuals.