

IMDB Movie Analysis

Kushagra Singh

March 2023

Description:

We have here dataset for IMDB Movies. I have framed the problem. For this task, I have define a problem I want to shed light on.

I did this by asking 'What?' i.e. What is the problem? Eg. What do I see happening? What is my hypothesis for the cause of the problem? What is the impact of the problem on stakeholders? What is the impact of the problem not being solved?

Answering these questions helped define problem I was trying to solve and will allowed me to find the right data to solve it.

Now next step I cleaned the data as necessary, and use your Data Analysis skills to explore the data set and derive insights.

Using Advanced Excel and Statistical concepts we have found out movies with the highest profit, top movies as per imdb rating, top directors, most popular genres, top foreign language films and more.

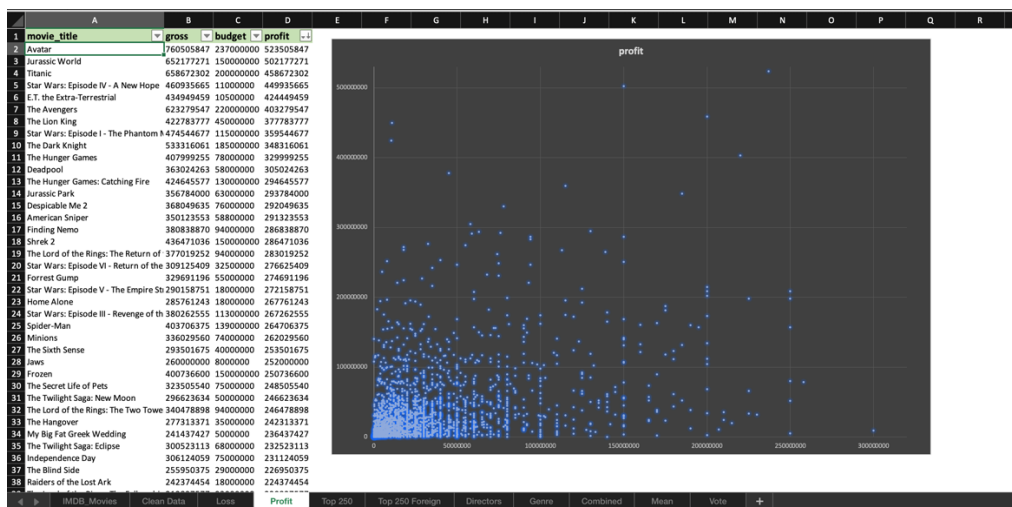
Data Cleaning:

Before analysing anything in data I have performed Data Cleaning.

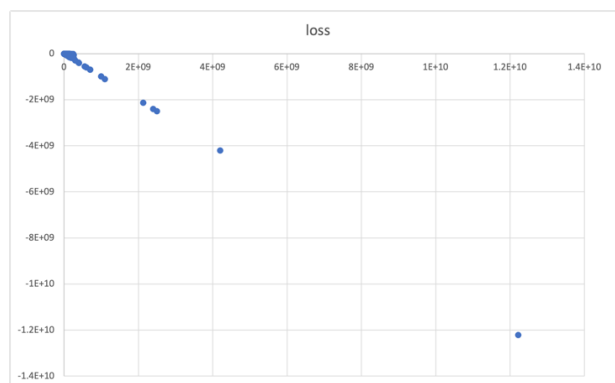
Where I performed three things

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

Q1. Which movie had the highest profit?



Avatar movie made the highest profit



Some observed loss too removing some outlier values gives us better idea



Here we have removed loss above 100000000

Interesting thing is movie Happiness made no profit no loss.

Approach: Created a new column called profit which contains the difference of the two columns: gross and budget. Sorted the column using the profit column as reference. Plotted profit (y-axis) vs budget (x- axis) and observed the outliers using the appropriate chart type.

Q2. What are the top 250 IMDB movies?

1	movie_title	num_voted_users	imdb_score	rank
2	The Shawshank Redemption	1689764	9.3	1
3	The Godfather	1155770	9.2	2
4	The Dark Knight	1676169	9	3
5	The Godfather: Part II	790926	9	4
6	Pulp Fiction	1324680	8.9	5
7	The Lord of the Rings: The Return of the King	1215718	8.9	6
8	Schindler's List	865020	8.9	7
9	The Good, the Bad and the Ugly	503509	8.9	8
10	Inception	1468200	8.8	9
11	Fight Club	1347461	8.8	10
12	Forrest Gump	1251222	8.8	11
13	The Lord of the Rings: The Fellowship of the Ring	1238746	8.8	12
14	Star Wars: Episode V - The Empire Strikes Back	837759	8.8	13
15	The Matrix	1217752	8.7	14
16	The Lord of the Rings: The Two Towers	1100446	8.7	15
17	Star Wars: Episode IV - A New Hope	911097	8.7	16
18	Goodfellas	728685	8.7	17
19	One Flew Over the Cuckoo's Nest	680041	8.7	18
20	City of God	533200	8.7	19
21	Seven Samurai	229012	8.7	20
22	Se7en	1023511	8.6	21
23	Interstellar	928227	8.6	22
24	The Silence of the Lambs	887467	8.6	23
25	Saving Private Ryan	881236	8.6	24
26	American History X	782437	8.6	25
27	The Usual Suspects	740918	8.6	26
28	Spirited Away	417971	8.6	27
29	Modern Times	143086	8.6	28
30	The Dark Knight Rises	1144337	8.5	29
31	Gladiator	982637	8.5	30
32	Django Unchained	955174	8.5	31
33	The Departed	873649	8.5	32
34	Memento	845580	8.5	33
35	The Prestige	844052	8.5	34
36	The Green Mile	782610	8.5	35
37	Terminator 2: Judgment Day	744891	8.5	36
38	Back to the Future	732212	8.5	37
39	Raiders of the Lost Ark	661017	8.5	38
40	The Lion King	644348	8.5	39
41	Alien	563827	8.5	40

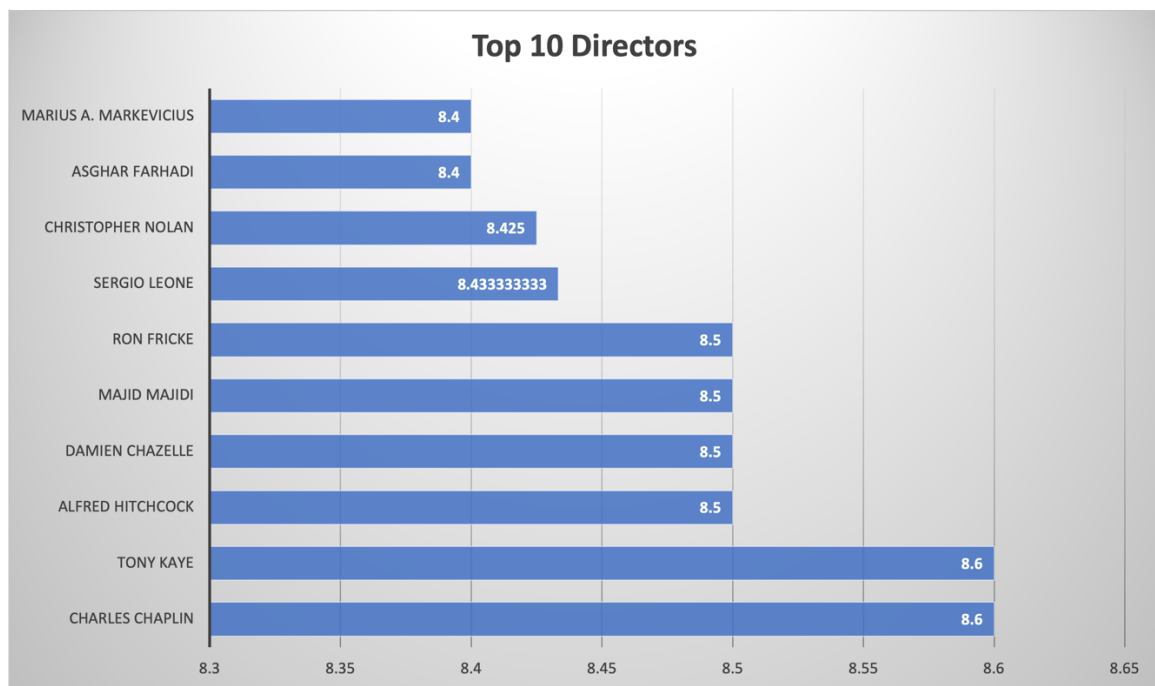
Approach: Created a new column IMDb_Top_250 and stored the top 250 movies with the highest IMDb Rating (corresponding to the column: imdb_score). Also made sure that for all of these movies, the num_voted_users is greater than 25,000. Also added a Rank column containing the values 1 to 250 indicating the ranks of the corresponding films.

Q3. Top 250 IMDB movies not in English?

	A	B	C	D	E
1	Top_Foreign_Lang_Film	num_voted_users	language	imdb_score	rank
2	The Good, the Bad and the Ugly	503509	Italian	8.9	1
3	City of God	533200	Portuguese	8.7	2
4	Seven Samurai	229012	Japanese	8.7	3
5	Spirited Away	417971	Japanese	8.6	4
6	The Lives of Others	259379	German	8.5	5
7	Children of Heaven	27882	Persian	8.5	6
8	Amélie	534262	French	8.4	7
9	Oldboy	356181	Korean	8.4	8
10	Princess Mononoke	221552	Japanese	8.4	9
11	Das Boot	168203	German	8.4	10
12	A Separation	151812	Persian	8.4	11
13	Baahubali: The Beginning	62756	Telugu	8.4	12
14	Downfall	248354	German	8.3	13
15	The Hunt	170155	Danish	8.3	14
16	Metropolis	111841	German	8.3	15
17	Pan's Labyrinth	467234	Spanish	8.2	16
18	Howl's Moving Castle	214091	Japanese	8.2	17
19	The Secret in Their Eyes	131831	Spanish	8.2	18
20	Incendies	80429	French	8.2	19
21	Amores Perros	173551	Spanish	8.1	20
22	Akira	106160	Japanese	8.1	21
23	Elite Squad	81644	Portuguese	8.1	22
24	The Celebration	65951	Danish	8.1	23
25	The Sea Inside	64556	Spanish	8.1	24
26	Tae Guk Gi: The Brotherhood of War	31943	Korean	8.1	25
27	A Fistful of Dollars	147566	Italian	8	26
28	Persepolis	70194	French	8	27
29	My Name Is Khan	69759	Hindi	8	28
30	Waltz with Bashir	46107	Hebrew	8	29
31	Central Station	28951	Portuguese	8	30
32	Crouching Tiger, Hidden Dragon	217740	Mandarin	7.9	31
33	Hero	149414	Mandarin	7.9	32
34	Letters from Iwo Jima	132149	Japanese	7.9	33
35	Amour	70382	French	7.9	34
36	4 Months, 3 Weeks and 2 Days	44763	Romanian	7.9	35
37	The Chorus	44151	French	7.9	36
38	Nine Queens	38215	Spanish	7.9	37
39	Veer-Zaara	34449	Hindi	7.9	38
40	Apocalypse	236000	Maya	7.8	39
41	Run Lola Run	361471	German	7.8	40

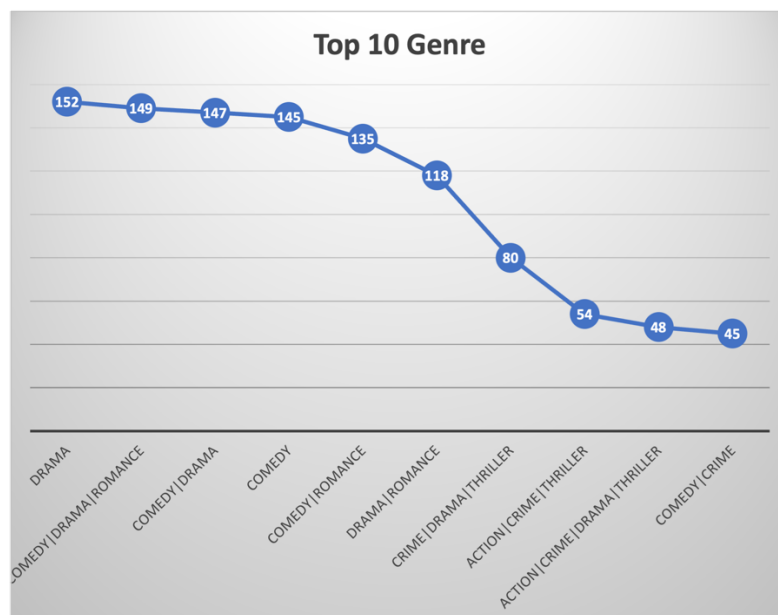
Approach: Extracted all the movies in the IMDb_Top_250 column which are not in the English language and stored them in a new column named Top_Foreign_Lang_Film.

Q3. Top 10 Directors as per imdb_rating?

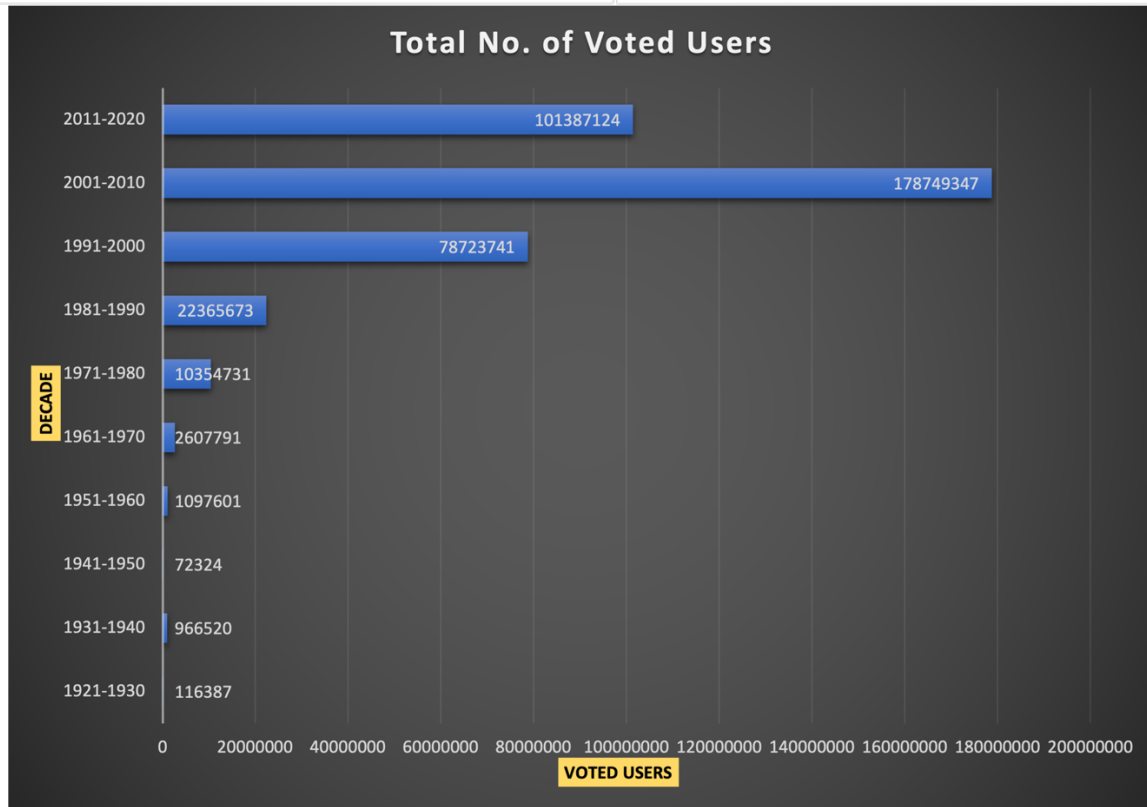
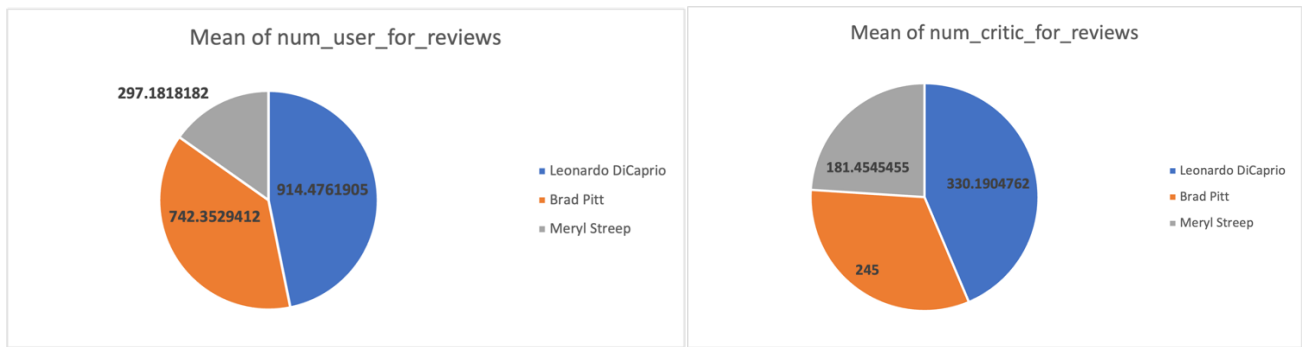


Approach: Grouped the column using the director_name column. Found out the top 10 directors for whom the mean of imdb_score is the highest and stored them in a new column top10director. In case of a tie in IMDb score between two directors, sorted them alphabetically. (graph lowest to highest)

Q4. Popular Genres?



Q5. Find the critic-favorite and audience-favorite actors?



Approach: Created three new columns for 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. Used only the actor_1_name column for extraction.

Appended the rows of all these columns and stored them in a new column named Combined. Grouped the combined column using the actor_1_name column.

Found the mean of the num_critic_for_reviews and num_users_for_review and identified the actors which have the highest mean.

Observed the change in number of voted users over decades using a bar chart. Created a column called decade which represents the decade to which every movie belongs

Data set used: [https://github.com/ks127d/Data-Analytics-Project/blob/dc7089f1f709c1b8e38250545efbb35a6312e51a/IMDB%20Movie%20Analysis/~\\$IMDB_Movies.xlsx](https://github.com/ks127d/Data-Analytics-Project/blob/dc7089f1f709c1b8e38250545efbb35a6312e51a/IMDB%20Movie%20Analysis/~$IMDB_Movies.xlsx)

(contains analysis too)