

# IMDB Movie Analysis

## Final Project-1

### Description:

For the Final Project, we are having a 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 we frame the problem i.e. What is the problem?

Using these questions which guide our thinking:

- 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?

Answering these questions will help you define a problem we are trying to solve and will allow you to find the right data to solve it.

Once we have defined a problem, clean the data as necessary, and use our Data Analysis skills to explore the data set and derive insights.

Make sure to use 5 Whys Analysis in your analysis and using this to create a report which conveys a data story.

Once we have framed the problem and gathered initial insights from the data, we can ask the following questions as you dig deeper into your analysis.

- What do you see happening?
- What are the specific symptoms of the problem?
- What is your hypothesis for the cause of the problem?

### Five 'Whys' approach

Once you have the problem better defined, you can use 5 Whys technique to determine

its root cause by repeatedly asking the question “Why”.

It's also called the Root Cause Analysis, developed by Sakichi Toyoda, founder of Toyota Industries. Here's an example of how this technique could be used to figure out the cause of the following problem: A business went over budget on a recent project.

Q: “Why did we go over budget on our project?”

A: It took much longer than we expected to complete.

Q: “Why did it take longer than expected to complete?”

A: We had to redesign several elements of the product.

Q: “Why did we have to redesign elements of the product?”

A: Features of the product were confusing to use.

Q: “Why were the features of the product confusing to use?”

A: We made incorrect assumptions about what users wanted.

Q: “Why did we make incorrect assumptions about what users wanted?”

A: Our user experience research team didn't ask effective questions.

As you see above, what looked like a budgeting problem turned out to be a problem with the user experience team not working effectively.

While asking Why is easy, what we're interested in is the answer. Each time you answer why the next time gets more difficult as you must think deeper behind the reasons for this. As you ask why, you may find that you have multiple answers for the same question.

**A. Cleaning the data:** This is one of the most important step to perform before moving forward with the analysis. Using our knowledge learned till now to do this like Dropping columns, removing null values, etc.  
To Clean the data.

## **Handling Duplicate Values**

- Press Ctrl+A to select the entire data
- Goto the data section
- Click on remove duplicates button
- Select all data and Select My data has headers in the pop-up menu
- Press ok to remove the duplicates value
- My data set had 44 duplicate rows , after removal 4999 unique rows exists.

## **Delete Empty Cells**

- Select the entire dataset by pressing Ctrl+A
- Select the 'find and select button' from home section.
- Then select go to special/ press F5 to access goto dialog box then click special
- Click blanks
- Click ok

Now the empty cells will get selected and we can modify it.

- To delete the empty cells
- Right click on one of the selected cell
- Click on delete option on the pop-up menu
- Now select delete entire row
- This will delete rows with empty cells

After deleting the empty cell rows we are having 3724 rows.

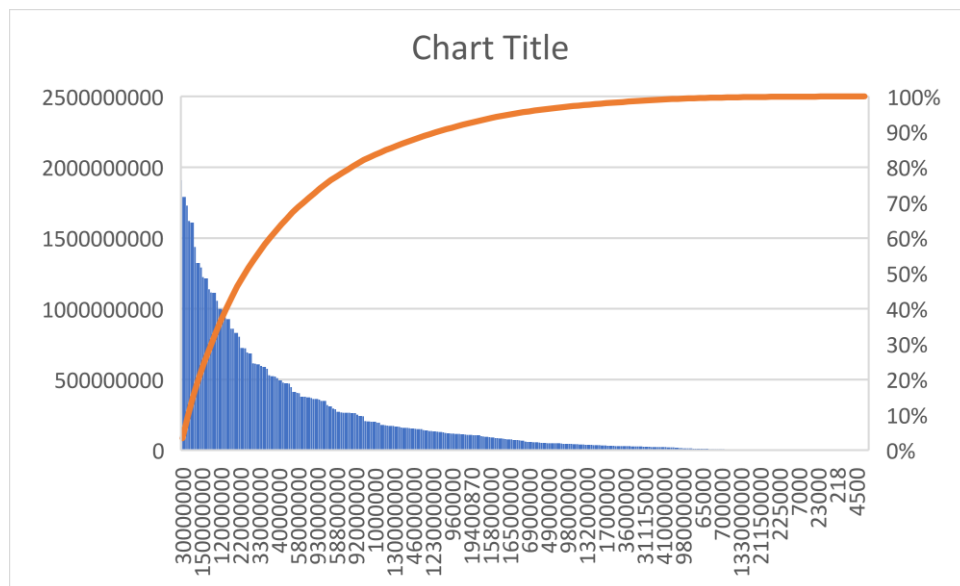
## **To Check Whether Empty Cells Exists**

- Select the entire column
- Press Ctrl+down
- If the pointer moves to the bottom last of the same column the empty cells does not exist.
- If it points to any empty cell in between again follow the process of deleting the empty rows
- I performed deleting the cells twice
- Then I checked until there is no empty cells by following these steps.

## **Alignment Of Data**

- To align the entire column
- Click on to the select column
- Now in the home section select left align option.
- Performing this operation will let us to align the data properly.

- B. 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.



- 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.

movie_title	imdb_score	RANK	num_voted_users
Avatar	9.3	1	886204
Pirates of the Caribbean: At World's End	9.2	2	471220
Spectre	9	3	275868
The Dark Knight Rises	9	4	1144337
John Carter	8.9	5	212204
Spider-Man 3	8.9	6	383056
Tangled	8.9	7	294810
Avengers: Age of Ultron	8.9	8	462669
Harry Potter and the Half-Blood Prince	8.8	9	321795
Batman v Superman: Dawn of Justice	8.8	10	371639
Superman Returns	8.8	11	240396
Quantum of Solace	8.8	12	330784
Pirates of the Caribbean: Dead Man's Chest	8.8	13	522040

The Lone Ranger	8.7	14	181792
Man of Steel	8.7	15	548573
The Chronicles of Narnia: Prince Caspian	8.7	16	149922
The Avengers	8.7	17	995415
Pirates of the Caribbean: On Stranger Tides	8.7	18	370704
Men in Black 3	8.7	19	268154
The Hobbit: The Battle of the Five Armies	8.7	20	354228
The Amazing Spider-Man	8.6	21	451803
Robin Hood	8.6	22	211765
The Hobbit: The Desolation of Smaug	8.6	23	483540
The Golden Compass	8.6	24	149019
King Kong	8.6	25	316018
Titanic	8.6	26	793059
Captain America: Civil War	8.6	27	272670
Battleship	8.6	28	202382
Jurassic World	8.5	29	418214
Skyfall	8.5	30	522030
Spider-Man 2	8.5	31	411164
Iron Man 3	8.5	32	557489
Alice in Wonderland	8.5	33	306320
X-Men: The Last Stand	8.5	34	383427
Monsters University	8.5	35	235025
Transformers: Revenge of the Fallen	8.5	36	323207
Transformers: Age of Extinction	8.5	37	242420
Oz the Great and Powerful	8.5	38	175409
The Amazing Spider-Man 2	8.5	39	321227
TRON: Legacy	8.5	40	264183
Cars 2	8.5	41	101178
Green Lantern	8.5	42	223393
Toy Story 3	8.5	43	544884
Terminator Salvation	8.5	44	286095
Furious 7	8.5	45	278232
World War Z	8.5	46	465019
X-Men: Days of Future Past	8.5	47	514125
Star Trek Into Darkness	8.4	48	395573
Jack the Giant Slayer	8.4	49	106416
The Great Gatsby	8.4	50	362912
Prince of Persia: The Sands of Time	8.4	51	222403
Pacific Rim	8.4	52	381148
Transformers: Dark of the Moon	8.4	53	326180
Indiana Jones and the Kingdom of the Crystal Skull	8.4	54	333847
Brave	8.4	55	273556
Star Trek Beyond	8.4	56	53607
WALL-E	8.4	57	718837
Rush Hour 3	8.4	58	121084
2012	8.4	59	283418
A Christmas Carol	8.4	60	72809

Jupiter Ascending	8.4	61	139593
The Legend of Tarzan	8.3	62	42372
The Chronicles of Narnia: The Lion, the Witch and the Wardrobe	8.3	63	286506
X-Men: Apocalypse	8.3	64	148379
The Dark Knight	8.3	65	1676169
Up	8.3	66	665575
Monsters vs. Aliens	8.3	67	114553
Iron Man	8.3	68	696338
Hugo	8.3	69	245333
Wild Wild West	8.3	70	129601
The Mummy: Tomb of the Dragon Emperor	8.3	71	117927
Suicide Squad	8.3	72	118992
Evan Almighty	8.3	73	115099
Edge of Tomorrow	8.3	74	431620
Waterworld	8.3	75	144337
G.I. Joe: The Rise of Cobra	8.3	76	174578
Inside Out	8.3	77	345198
The Jungle Book	8.3	78	106072
Iron Man 2	8.3	79	522371
Snow White and the Huntsman	8.3	80	228554
Maleficent	8.3	81	252257
Dawn of the Planet of the Apes	8.3	82	317542
47 Ronin	8.3	83	116994
Captain America: The Winter Soldier	8.3	84	496749
Shrek Forever After	8.3	85	138661
Tomorrowland	8.2	86	128306
Big Hero 6	8.2	87	279093
Wreck-It Ralph	8.2	88	272534
The Polar Express	8.2	89	120798
Independence Day: Resurgence	8.2	90	58137
How to Train Your Dragon	8.2	91	485430
Terminator 3: Rise of the Machines	8.2	92	305340
Guardians of the Galaxy	8.2	93	682155
Interstellar	8.2	94	928227
Inception	8.2	95	1468200
The Fast and the Furious	8.2	96	272223
The Curious Case of Benjamin Button	8.2	97	459346
X-Men: First Class	8.2	98	518537
The Hunger Games: Mockingjay - Part 2	8.2	99	166137
The Sorcerer's Apprentice	8.2	100	124185
Poseidon	8.2	101	82380
Warcraft	8.2	103	211971
Terminator Genisys	8.2	104	111609
The Chronicles of Narnia: The Voyage of the Dawn Treader	8.2	105	188457
Pearl Harbor	8.2	106	106446
Transformers	8.2	107	254111

AlexanderÂ	8.2	108	513158
Harry Potter and the Order of the PhoenixÂ	8.2	109	138863
Harry Potter and the Goblet of FireÂ	8.1	110	355137
HancockÂ	8.1	111	385670
I Am LegendÂ	8.1	112	343648
Charlie and the Chocolate FactoryÂ	8.1	113	530870
RatatouilleÂ	8.1	114	320284
Batman BeginsÂ	8.1	115	473887
Madagascar: Escape 2 AfricaÂ	8.1	116	980946
Night at the Museum: Battle of the SmithsonianÂ	8.1	117	146019
X-Men Origins: WolverineÂ	8.1	118	130272
The Matrix RevolutionsÂ	8.1	119	361924
FrozenÂ	8.1	120	364948
The Matrix ReloadedÂ	8.1	121	421658
Thor: The Dark WorldÂ	8.1	122	421818
Mad Max: Fury RoadÂ	8.1	123	414070
Angels & DemonsÂ	8.1	124	552503
ThorÂ	8.1	125	207839
BoltÂ	8.1	126	536314
G-ForceÂ	8.1	127	146766
Wrath of the TitansÂ	8.1	128	33042
Dark ShadowsÂ	8.1	129	152826
Mission: Impossible - Rogue NationÂ	8.1	130	199039
The WolfmanÂ	8.1	131	232187
Bee MovieÂ	8.1	132	89442
Kung Fu Panda 2Â	8.1	133	105902
The Last AirbenderÂ	8.1	134	182718
Mission: Impossible IIIÂ	8.1	135	118951
White House DownÂ	8.1	136	256695
Flushed AwayÂ	8.1	137	164238
Mr. Peabody & ShermanÂ	8.1	139	85086
TroyÂ	8.1	140	39956
Madagascar 3: Europe's Most WantedÂ	8.1	141	47900
Die Another DayÂ	8.1	142	381672
GhostbustersÂ	8.1	143	119213
ArmageddonÂ	8.1	144	169914
Men in Black IIÂ	8.1	145	69757
BeowulfÂ	8.1	146	322395
Kung Fu Panda 3Â	8.1	147	270207
Mission: Impossible - Ghost ProtocolÂ	8.1	148	142440
Rise of the GuardiansÂ	8.1	149	64322
Fun with Dick and JaneÂ	8.1	150	365104
The Last SamuraiÂ	8.1	151	123553
Exodus: Gods and KingsÂ	8.1	152	110788
Star TrekÂ	8.1	153	317166
Spider-ManÂ	8.1	154	128682
How to Train Your Dragon 2Â	8.1	155	504419

Gods of Egypt	8.1	156	544665
Stealth	8	157	221128
Watchmen	8	158	51892
Lethal Weapon 4	8	159	45455
Hulk	8	160	392474
G.I. Joe: Retaliation	8	161	127497
Sahara	8	162	212106
Final Fantasy: The Spirits Within	8	163	146352
Captain America: The First Avenger	8	164	77673
The World Is Not Enough	8	165	72259
Master and Commander: The Far Side of the World	8	166	508818
The Twilight Saga: Breaking Dawn - Part 2	8	167	157519
Happy Feet 2	8	168	168207
The Incredible Hulk	8	169	185394
The Revenant	8	170	32399
Turbo	8	171	326286
Penguins of Madagascar	8	173	406020
The Bourne Ultimatum	8	174	62424
Kung Fu Panda	8	175	183208
Ant-Man	8	176	60230
The Hunger Games: Catching Fire	8	177	491077
Home	8	178	307029
War of the Worlds	8	179	313866
Bad Boys II	8	180	498397
Puss in Boots	8	181	70121
Salt	8	182	334345
Noah	8	183	178126
The Adventures of Tintin	8	184	114287
Harry Potter and the Prisoner of Azkaban	8	185	245621
Australia	8	186	200022
After Earth	8	187	177383
Dinosaur	8	188	382255
Night at the Museum: Secret of the Tomb	8	189	102338
Megamind	8	190	158720
Harry Potter and the Sorcerer's Stone	8	191	38438
R.I.P.D.	8	192	67223
Pirates of the Caribbean: The Curse of the Black Pearl	8	193	172754
The Hunger Games: Mockingjay - Part 1	8	194	444683
The Da Vinci Code	8	195	91640
Rio 2	8	196	809474
X-Men 2	8	197	305008
Fast Five	8	198	314253
Sherlock Holmes: A Game of Shadows	8	199	58498
Clash of the Titans	8	200	405973
Total Recall	8	201	284792
The 13th Warrior	8	202	338635
The Bourne Legacy	8	203	229679



Batman & Robin	8	204	240241
How the Grinch Stole Christmas	8	205	101411
The Day After Tomorrow	8	206	229823
Mission: Impossible II	8	207	189855
The Perfect Storm	8	208	141414
Fantastic 4: Rise of the Silver Surfer	8	209	333248
Life of Pi	8	210	242188
Ghost Rider	7.9	211	133076
Jason Bourne	7.9	212	213275
Charlie's Angels: Full Throttle	7.9	213	440084
Prometheus	7.9	214	182661
Stuart Little 2	7.9	215	40123
Elysium	7.9	216	100821
The Chronicles of Riddick	7.9	217	456260
RoboCop	7.9	218	36471
Speed Racer	7.9	219	338087
How Do You Know	7.9	220	183909
Knight and Day	7.9	221	182899
Oblivion	7.9	222	57873
Star Wars: Episode III - Revenge of the Sith	7.9	223	35066
Star Wars: Episode II - Attack of the Clones	7.9	224	148280
Monsters, Inc.	7.9	225	387436
The Wolverine	7.9	226	520104
Star Wars: Episode I - The Phantom Menace	7.9	227	464310
The Croods	7.9	228	585659
Windtalkers	7.9	229	328067
The Huntsman: Winter's War	7.9	230	534658
Teenage Mutant Ninja Turtles	7.9	231	150618
Gravity	7.9	232	55994
Dante's Peak	7.9	233	37750
Fantastic Four	7.9	234	167085
Night at the Museum	7.9	235	582917
San Andreas	7.9	236	62271
Tomorrow Never Dies	7.9	237	110486
The Patriot	7.9	238	234480
Ocean's Twelve	7.9	239	147497
Mr. & Mrs. Smith	7.9	240	149680
Insurgent	7.9	241	207613
The Aviator	7.9	242	284852
Gulliver's Travels	7.9	243	348861
The Green Hornet	7.9	244	154621
300: Rise of an Empire	7.9	245	264318
The Smurfs	7.9	246	53160
Allegiant	7.9	247	136019
Real Steel	7.9	248	225273
The Smurfs 2	7.9	249	66593
Ender's Game	7.9	250	44296

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!

Rank	language	movie_title
1	Mandarin	The Flowers of War
2	Aboriginal	The Interpreter
3	Spanish	The Legend of Zorro
4	French	Oceans
5	Mandarin	Dragon Blade
6	Filipino	The Great Raid
7	French	A Very Long Engagement
		Curse of the Golden
8	Mandarin	Flower
9	Mandarin	Hero
10	French	Micmacs
11	Maya	Apocalypto
12	French	Amélie
13	Mandarin	The Warlords
14	Kazakh	Nomad: The Warrior
15	Mandarin	Red Cliff
16	Mandarin	The Grandmaster
17	Cantonese	Ip Man 3
18	language	movie_title
19	Mandarin	The Flowers of War
20	Aboriginal	The Interpreter
21	Spanish	The Legend of Zorro
22	French	Oceans
23	Mandarin	Dragon Blade
24	Filipino	The Great Raid
25	French	A Very Long Engagement
		Curse of the Golden
26	Mandarin	Flower
27	Mandarin	Hero
28	French	Micmacs
29	Maya	Apocalypto
30	French	Amélie
31	Mandarin	The Warlords
32	Kazakh	Nomad: The Warrior
33	Mandarin	Red Cliff
34	Mandarin	The Grandmaster
35	Cantonese	Ip Man 3
36	language	movie_title
37	Mandarin	The Flowers of War
38	Aboriginal	The Interpreter

39	Spanish	The Legend of Zorro
40	French	Oceans
41	Mandarin	Dragon Blade
42	language	movie_title
43	Mandarin	The Flowers of War

**D. Best Directors:** Group the column using the director\_name column.

Finding 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.

director_name	imdb_score
Christopher Nolan	9
Christopher Nolan	8.8
David Fincher	8.8
Francis Ford Coppola	9
Francis Ford Coppola	9.2
Frank Darabont	9.3
Irvin Kershner	8.8
Peter Jackson	8.8
Peter Jackson	8.9
Quentin Tarantino	8.9
Robert Zemeckis	8.8
Sergio Leone	8.9
Steven Spielberg	8.9

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

#### genres

Action|Crime|Drama|Thriller  
 Action|Adventure|Sci-Fi|Thriller  
 Action|Adventure|Drama|Fantasy  
 Action|Adventure|Drama|Fantasy  
 Action|Sci-Fi  
 Drama  
 Comedy|Drama

Crime|Drama|Sci-Fi  
Crime|Drama  
Crime|Drama|Thriller

- F. **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.

Meryl_Streep	Leo_Caprio	Brad_Pitt
It's ComplicatedÂ	TitanicÂ	The Curious Case of Benjamin ButtonÂ
The River WildÂ	The Great GatsbyÂ	TroyÂ
Julie & JuliaÂ	InceptionÂ	Ocean's TwelveÂ
The Devil Wears PradaÂ	The RevenantÂ	Mr. & Mrs. SmithÂ
Lions for LambsÂ	The AviatorÂ	Spy GameÂ
Out of AfricaÂ	Django UnchainedÂ	Ocean's ElevenÂ
Hope SpringsÂ	Blood DiamondÂ	FuryÂ
One True ThingÂ	The Wolf of Wall StreetÂ	Seven Years in TibetÂ
The HoursÂ	Gangs of New YorkÂ	Fight ClubÂ
The Iron LadyÂ	The DepartedÂ	Sinbad: Legend of the Seven SeasÂ
A Prairie Home CompanionÂ	Shutter IslandÂ	Interview with the Vampire: The Vampire ChroniclesÂ
	Body of LiesÂ	The Tree of LifeÂ
	Catch Me If You CanÂ	The Assassination of Jesse James by the Coward Robert FordÂ
	The BeachÂ	BabelÂ
	Revolutionary RoadÂ	By the SeaÂ
	The Man in the Iron MaskÂ	Killing Them SoftlyÂ
	J. EdgarÂ	True RomanceÂ
	The Quick and the DeadÂ	
	Marvin's RoomÂ	
	Romeo + JulietÂ	
	The Great GatsbyÂ	

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.

ACTOR	COMBINED
Meryl_Streep	It's ComplicatedÂ
Meryl_Streep	The River WildÂ
Meryl_Streep	Julie & JuliaÂ
Meryl_Streep	The Devil Wears PradaÂ
Meryl_Streep	Lions for LambsÂ
Meryl_Streep	Out of AfricaÂ
Meryl_Streep	Hope SpringsÂ
Meryl_Streep	One True ThingÂ
Meryl_Streep	The HoursÂ
Meryl_Streep	The Iron LadyÂ
Meryl_Streep	A Prairie Home CompanionÂ
Leo_Caprio	TitanicÂ
Leo_Caprio	The Great GatsbyÂ
Leo_Caprio	InceptionÂ
Leo_Caprio	The RevenantÂ
Leo_Caprio	The AviatorÂ
Leo_Caprio	Django UnchainedÂ
Leo_Caprio	Blood DiamondÂ
Leo_Caprio	The Wolf of Wall StreetÂ
Leo_Caprio	Gangs of New YorkÂ
Leo_Caprio	The DepartedÂ
Leo_Caprio	Shutter IslandÂ
Leo_Caprio	Body of LiesÂ
Leo_Caprio	Catch Me If You CanÂ
Leo_Caprio	The BeachÂ
Leo_Caprio	Revolutionary RoadÂ
Leo_Caprio	The Man in the Iron MaskÂ
Leo_Caprio	J. EdgarÂ
Leo_Caprio	The Quick and the DeadÂ
Leo_Caprio	Marvin's RoomÂ
Leo_Caprio	Romeo + JulietÂ
Leo_Caprio	The Great GatsbyÂ
Brad_Pitt	The Curious Case of Benjamin ButtonÂ
Brad_Pitt	TroyÂ
Brad_Pitt	Ocean's TwelveÂ
Brad_Pitt	Mr. & Mrs. SmithÂ
Brad_Pitt	Spy GameÂ
Brad_Pitt	Ocean's ElevenÂ
Brad_Pitt	FuryÂ
Brad_Pitt	Seven Years in TibetÂ

Brad_Pitt	Fight Club
Brad_Pitt	Sinbad: Legend of the Seven Seas
Brad_Pitt	Interview with the Vampire: The Vampire Chronicles
Brad_Pitt	The Tree of Life
	The Assassination of Jesse James by the Coward Robert
Brad_Pitt	Ford
Brad_Pitt	Babel
Brad_Pitt	By the Sea
Brad_Pitt	Killing Them Softly
Brad_Pitt	True Romance

Find the mean of the num\_critic\_for\_reviews and num\_users\_for\_review and identify the actors which have the highest mean.

#### **G. num\_critic\_for\_reviews**

MEAN MERYL STREEP - 181.4545

MEAN Leo\_Caprio - 330.1905

MEAN Brad\_Pitt - 245

#### **H. num\_users\_for\_review**

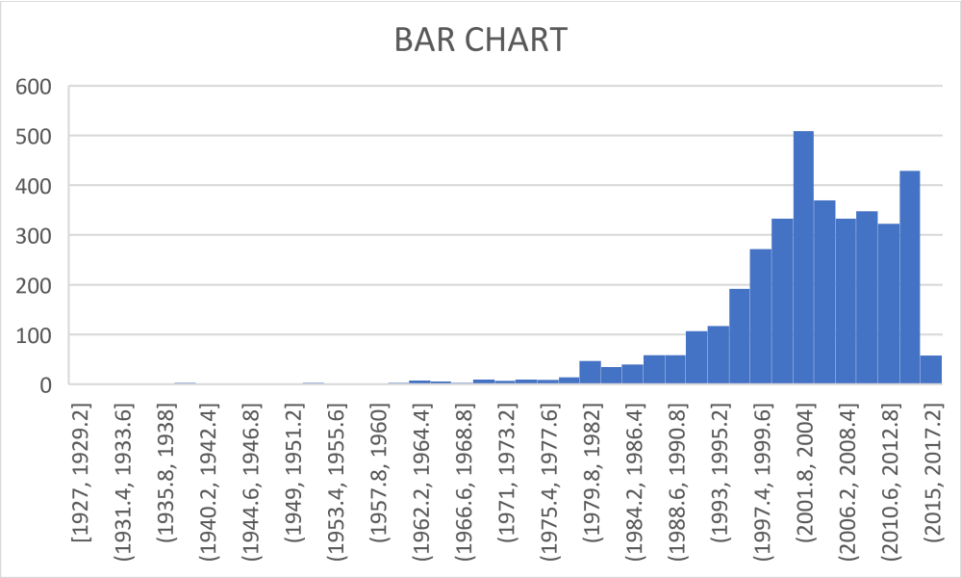
MEAN MERYL STREEP- 297.1818

MEAN Leo\_Caprio- 914.4762

MEAN Brad\_Pitt - 742.3529

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.



The critic-favorite and audience-favorite actor is Leo\_Caprio!