



Recommendation System for Movies

An analysis of Movie lens Data

Topic Overview

- Problem
- Dataset
- Analysis and Result
- Machine Learning
- Conclusion

Why we need Recommendation system?

- It has been estimated that there are approximately 500,000 movies currently in existence.
- This situation creates 'selection pressure' on individual movies, as not all movies are equally popular.
- The audience will not be spending their time watching every movie available to them, they will pick randomly something to watch.
- If a person picks a movie and does not enjoy it there will be no positive word-of-mouth.
- This doesn't necessarily mean the movie was bad, it might mean it was not interesting to that individual. This is where the recommendation system is helpful.
- Recommendation system helps the user find items of their interest and Helps the item provider to deliver their items to the right user.
- It increases revenues for business through increased consumption.

Dataset

- metadata for all 45,000 movies listed in the Full Movie Lens Dataset.
- The dataset consists of movies released on or before July 2017.
- This dataset also has files containing 26 million ratings from 270,000 users for all 45,000 movies.

```
RangeIndex: 45466 entries, 0 to 45465
Data columns (total 24 columns):
adult                45466 non-null object
belongs_to_collection  4494 non-null object
budget              45466 non-null object
genres              45466 non-null object
homepage            7782 non-null object
id                  45466 non-null object
imdb_id             45449 non-null object
original_language    45455 non-null object
original_title       45466 non-null object
overview            44512 non-null object
popularity           45461 non-null object
poster_path          45080 non-null object
production_companies  45463 non-null object
production_countries  45463 non-null object
release_date         45379 non-null object
revenue              45460 non-null float64
runtime              45203 non-null float64
spoken_languages     45460 non-null object
status               45379 non-null object
tagline              20412 non-null object
title                45460 non-null object
video                45460 non-null object
vote_average          45460 non-null float64
vote_count            45460 non-null float64
dtypes: float64(4), object(20)
```

Data Wrangling

- The dataset has a record of 45466 movies with 24 columns(features).Dropped columns with more than 50% of its total data with null values.
- Converted budget feature into a numeric variable and replaced all the non-numeric values with NaN.
- Extracted feature "Release Year" from "Release Date." "Release Year" is the year in which the movie was released.
- Calculated the Net Profit/Loss using features "Revenue" and "budget."
- Dropped features like adult, id, original_title, poster_path, video as does not provide useful information.
- Converted vote_count, and vote_average to float type from object.

Which Production companies make most money in movie business?

- Warner Bros is the highest-earning production company of all time, earning a staggering 63.5 billion dollars from close to 500 movies.
- Universal Pictures and Paramount Pictures are the second and the third highest-earning companies with 55 billion dollars and 48 billion dollars in revenue.
- As we are aware, Warner Bros and Universal Pictures are bigger studios compared to others on the list.
- Thus it would be more appropriate to look at the average revenue of studios. We will consider studios that have produced at least ten movies.

	Total Revenue	Average Revenue	Number Of Movies
Warner Bros.	6.352519e+10	1.293792e+08	491
Universal Pictures	5.525919e+10	1.193503e+08	463
Paramount Pictures	4.880819e+10	1.235650e+08	395
Twentieth Century Fox Film Corporation	4.768775e+10	1.398468e+08	341
Walt Disney Pictures	4.083727e+10	2.778046e+08	147
Columbia Pictures	3.227974e+10	1.367785e+08	236
New Line Cinema	2.217339e+10	1.119868e+08	198
Amblin Entertainment	1.734372e+10	2.550547e+08	68
DreamWorks SKG	1.547575e+10	1.984071e+08	78
Dune Entertainment	1.500379e+10	2.419966e+08	62

Which Production companies make most money in movie business?

- Pixar Animation Studios has produced the most successful movies, on average.
- Marvel Studios, with an average gross of 615 million dollars, comes in second.

	Total Revenue	Average Revenue	Number of Movies
Pixar Animation Studios	1.118853e+10	6.215852e+08	37
Marvel Studios	1.169964e+10	6.157703e+08	19
Heyday Films	7.920012e+09	6.092317e+08	13
WingNut Films	7.111004e+09	5.470003e+08	13
Revolution Sun Studios	8.120339e+09	5.413559e+08	15
Syncopy	5.359856e+09	5.359856e+08	10
Fuji Television Network	5.880444e+09	4.900370e+08	12
Blue Sky Studios	5.274028e+09	4.794570e+08	11
Walt Disney Animation Studios	6.053112e+09	4.656240e+08	13
Lucasfilm	9.898421e+09	4.499282e+08	22

Which movies are more popular?

- "Minions" is the most popular movie by the TMDB Popularity Score.
- Wonder Woman and Beauty and the Beast come in second and third respectively.

	title	popularity	year
30700	Minions	547.488298	2015.0
33356	Wonder Woman	294.337037	2017.0
42222	Beauty and the Beast	287.253654	2017.0
43644	Baby Driver	228.032744	2017.0
24455	Big Hero 6	213.849907	2014.0
26564	Deadpool	187.860492	2016.0
26566	Guardians of the Galaxy Vol. 2	185.330992	2017.0
14551	Avatar	185.070892	2009.0
24351	John Wick	183.870374	2014.0
23675	Gone Girl	154.801009	2014.0

Which movies have been most voted by TMDB voters?

- Inception and The Dark Knight, two critically acclaimed movies, are at the top of our chart.
- It is interesting to note that Christopher Nolan directed both movies.

	title	vote_count	year
15480	Inception	14075.0	2010.0
12481	The Dark Knight	12269.0	2008.0
14551	Avatar	12114.0	2009.0
17818	The Avengers	12000.0	2012.0
26564	Deadpool	11444.0	2016.0
22879	Interstellar	11187.0	2014.0
20051	Django Unchained	10297.0	2012.0
23753	Guardians of the Galaxy	10014.0	2014.0
2843	Fight Club	9678.0	1999.0
18244	The Hunger Games	9634.0	2012.0

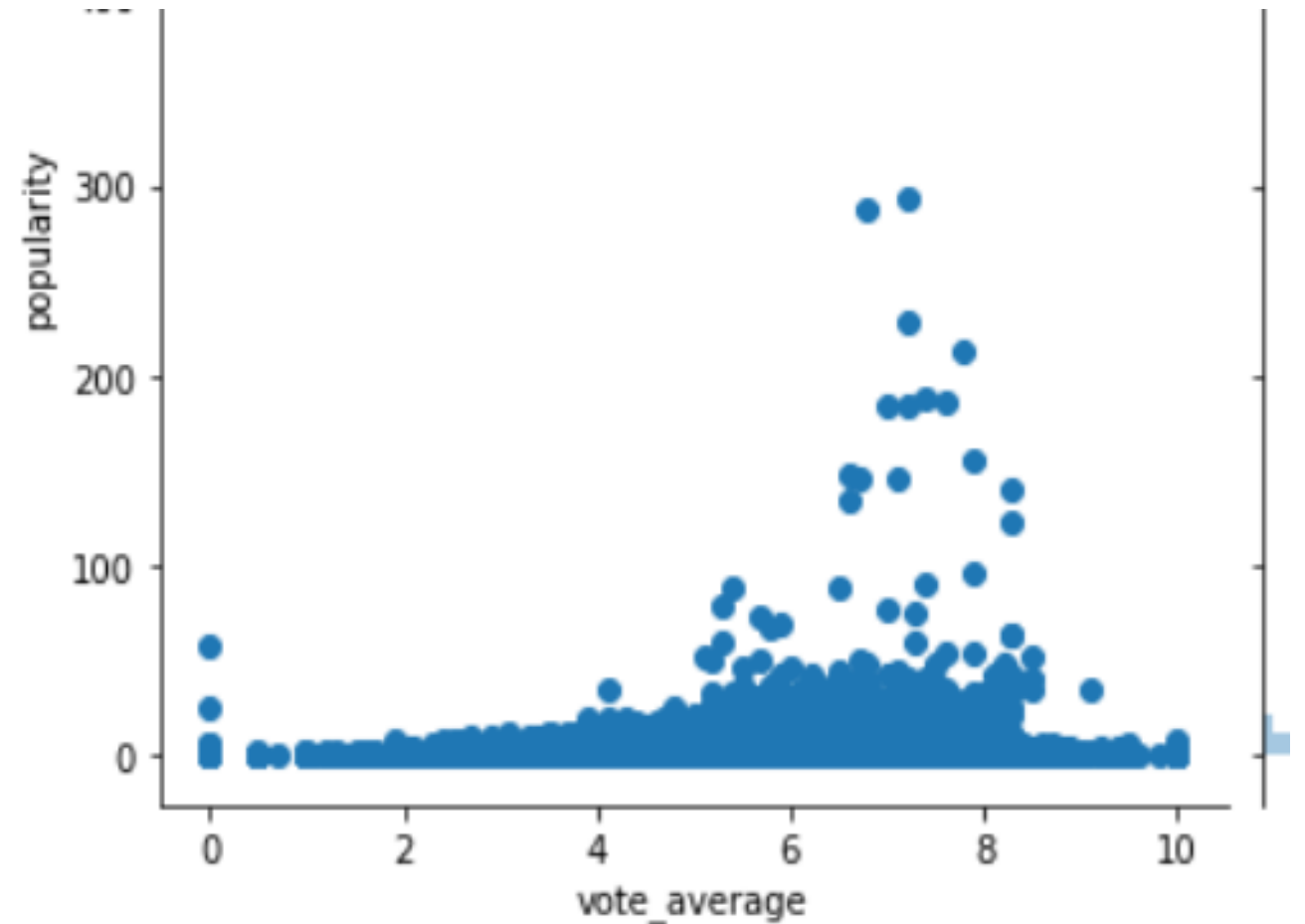
Which movies are most Critically Acclaimed?

- We will only consider those movies with more than 5000 votes (similar to IMDB's criteria of 5000 options in selecting its top 250).
- The Shawshank Redemption and The Godfather are the two most critically acclaimed movies in the TMDB Database.

	title	vote_average	vote_count	year
314	The Shawshank Redemption	8.5	8358.0	1994.0
834	The Godfather	8.5	6024.0	1972.0
292	Pulp Fiction	8.3	8670.0	1994.0
12481	The Dark Knight	8.3	12269.0	2008.0
2843	Fight Club	8.3	9678.0	1999.0
18465	The Intouchables	8.2	5410.0	2011.0
351	Forrest Gump	8.2	8147.0	1994.0
1154	The Empire Strikes Back	8.2	5998.0	1980.0
256	Star Wars	8.1	6778.0	1977.0
46	Se7en	8.1	5915.0	1995.0

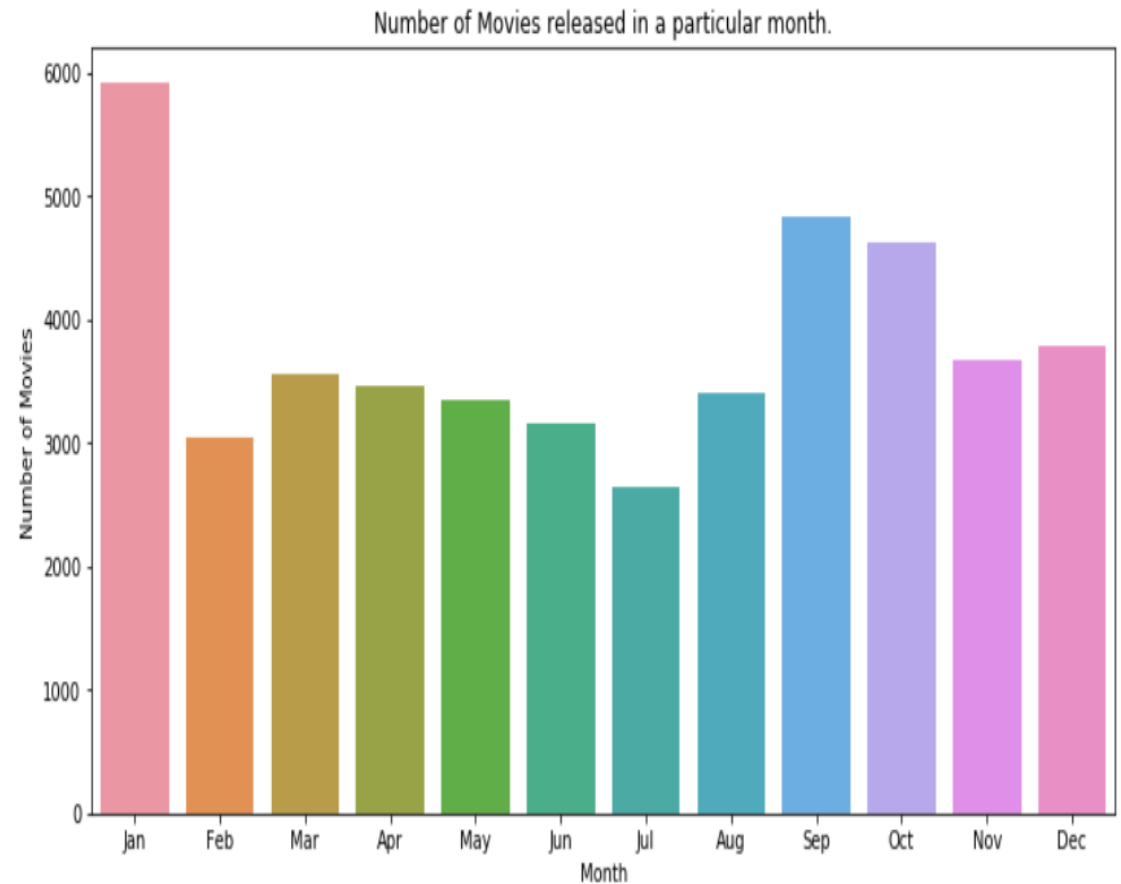
Do popularity and vote average share a tangible relationship?

- The Pearson Coefficient of the two quantities, is 0.154, suggesting no definite correlation.
- In other words, popularity and vote average are independent quantities.



Does Release Dates play a significant role in determining the revenue of a movie?

- We will gain insights about release dates in terms of months.
- Extracted the month for each movie from release date.
- January is the most popular month when it comes to movie releases.
- This is also known as the dump month in Hollywood circles when the dozen release subpar movies.



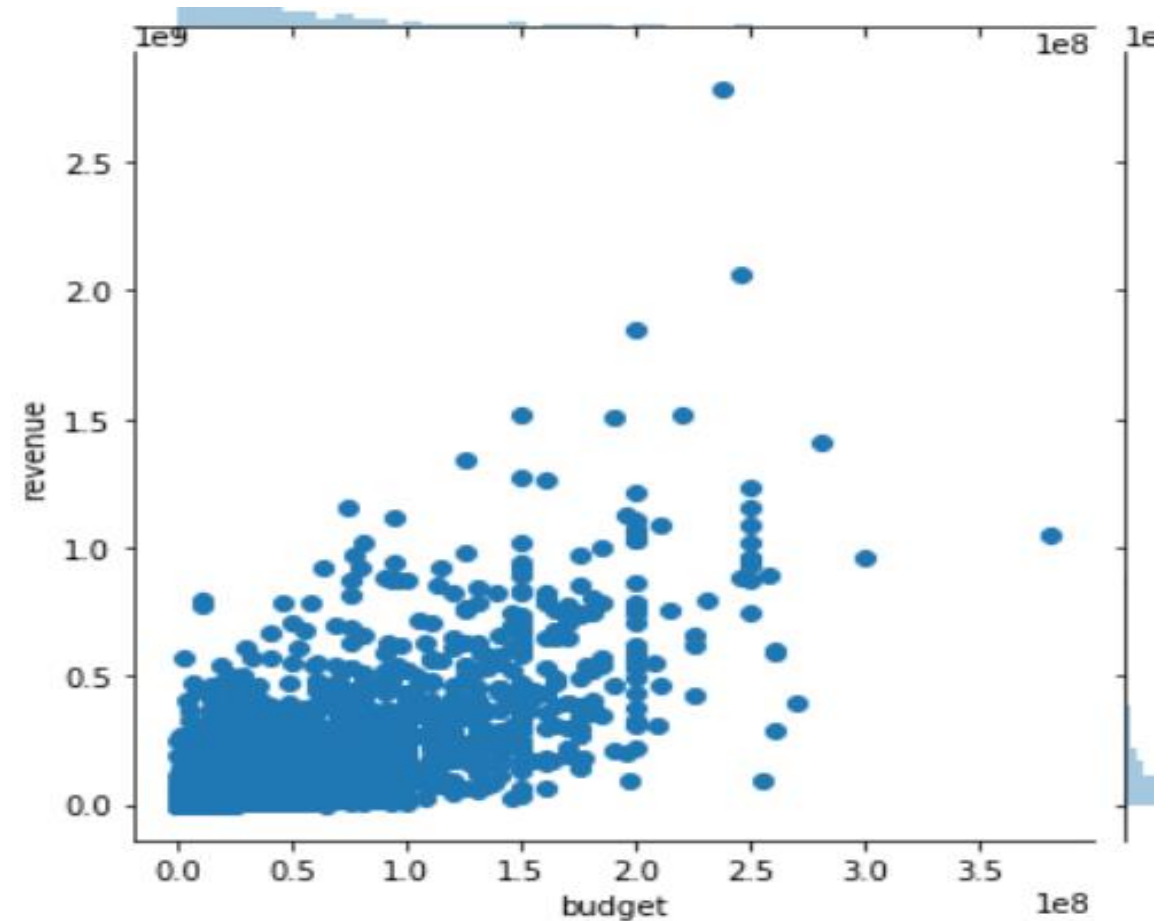
Which are the most expensive movies of all time?

- Two Pirates of the Caribbean films occupy the top spots in this list with a staggering budget of over 300 million dollars.
- All the top 10 most expensive movies made a profit on their investment except for The Lone Ranger
- “The Lone Ranger “ managed to recoup less than 35% of its investment, taking in a paltry 90 million dollars on a 255-million-dollar budget.

	title	budget	revenue	Net Profit/Loss	year
17124	Pirates of the Caribbean: On Stranger Tides	380000000.0	1.045714e+09	2.751878	2011.0
11827	Pirates of the Caribbean: At World's End	300000000.0	9.610000e+08	3.203333	2007.0
26558	Avengers: Age of Ultron	280000000.0	1.405404e+09	5.019299	2015.0
11067	Superman Returns	270000000.0	3.910812e+08	1.448449	2006.0
44842	Transformers: The Last Knight	260000000.0	6.049421e+08	2.326701	2017.0
16130	Tangled	260000000.0	5.917949e+08	2.276134	2010.0
18685	John Carter	260000000.0	2.841391e+08	1.092843	2012.0
11780	Spider-Man 3	258000000.0	8.908716e+08	3.452991	2007.0
21175	The Lone Ranger	255000000.0	8.928991e+07	0.350157	2013.0
22059	The Hobbit: The Desolation of Smaug	250000000.0	9.584000e+08	3.833600	2013.0

How strong a correlation does the budget hold with the revenue?

- A stronger correlation would directly imply more accurate forecasts.
- The scatterplot shows a positive correlation between budget and revenue.



Which are the Highest Grossing Films of All Time?

- The world of movies broke the 1-billion-dollar mark in 1997 with the release of Titanic.
- It took another 12 years to break the 2-billion-dollar mark with Avatar. James Cameron directed both these movies.
- The highest-grossing movie does not necessarily mean the movie made the highest profit of all.

	title	budget	revenue	year
14551	Avatar	237000000.0	2.787965e+09	2009.0
26555	Star Wars: The Force Awakens	245000000.0	2.068224e+09	2015.0
1639	Titanic	200000000.0	1.845034e+09	1997.0
17818	The Avengers	220000000.0	1.519558e+09	2012.0
25084	Jurassic World	150000000.0	1.513529e+09	2015.0
28830	Furious 7	190000000.0	1.506249e+09	2015.0
26558	Avengers: Age of Ultron	280000000.0	1.405404e+09	2015.0
17437	Harry Potter and the Deathly Hallows: Part 2	125000000.0	1.342000e+09	2011.0
22110	Frozen	150000000.0	1.274219e+09	2013.0
42222	Beauty and the Beast	160000000.0	1.262886e+09	2017.0

Which are the least and the most successful movies of all time?

- Let us check the least and the most successful movies of all time.
- we will only consider those movies which have a budget greater than 5 million dollars.
- E.T. the Extra-Terrestrial is the most successful movie.
- It is interesting to note that most of the successful movies in the top 10 list are released between 1965 - 1989.

	title	budget	revenue	Net Profit/Loss	year
1065	E.T. the Extra-Terrestrial	10500000.0	792965326.0	75.520507	1982.0
256	Star Wars	11000000.0	775398007.0	70.490728	1977.0
1338	Jaws	7000000.0	470654000.0	67.236286	1975.0
1888	The Exorcist	8000000.0	441306145.0	55.163268	1973.0
352	Four Weddings and a Funeral	6000000.0	254700832.0	42.450139	1994.0
834	The Godfather	6000000.0	245066411.0	40.844402	1972.0
4492	Look Who's Talking	7500000.0	296000000.0	39.466667	1989.0
24258	Annabelle	6500000.0	255273813.0	39.272894	2014.0
1056	Dirty Dancing	6000000.0	213954274.0	35.659046	1987.0
1006	The Sound of Music	8200000.0	286214286.0	34.904181	1965.0

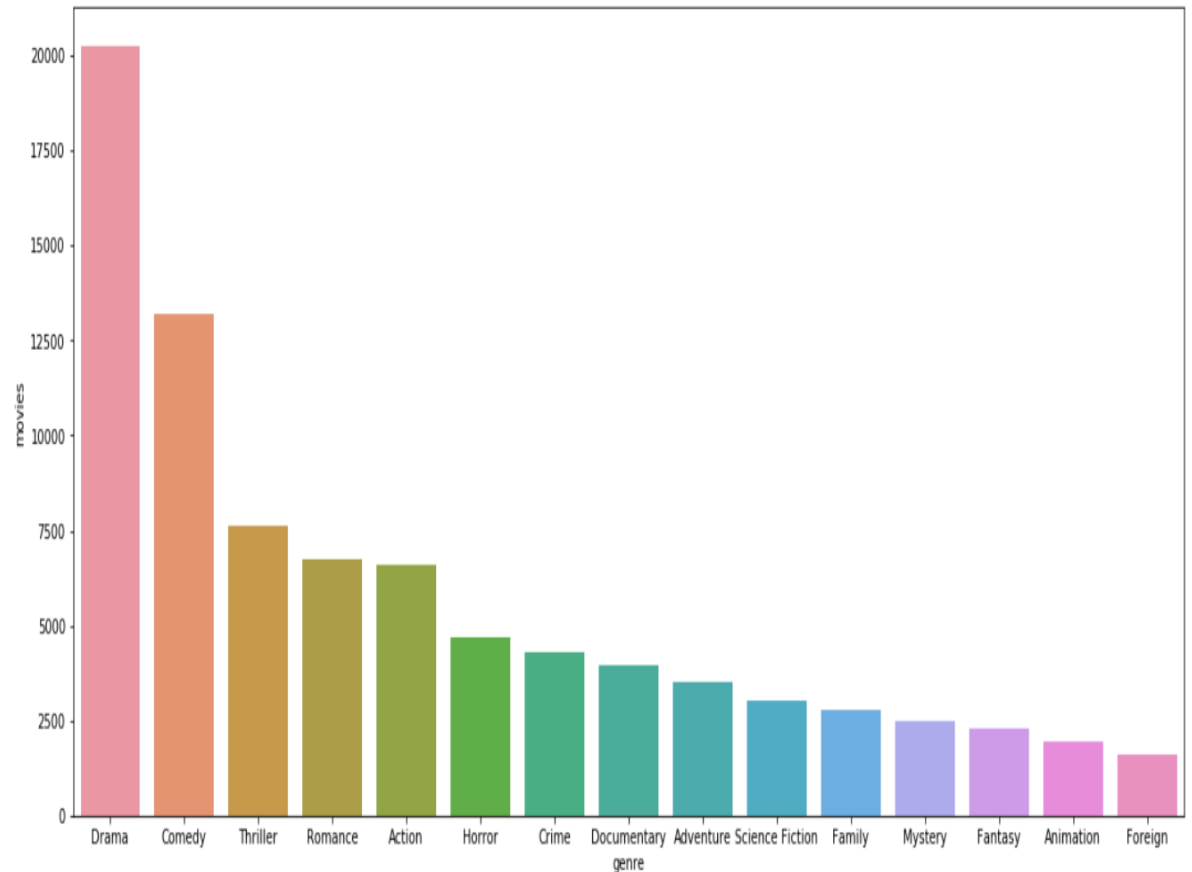
Which are the least and the most successful movies of all time?

- Chaos is the least successful movie.
- We can observe here that most of the movies listed in top 10 are released between 2000-2012
- So is it that older movies were more successful than newer ones?
- We cannot certainly say so as these figures have not been adjusted for inflation.

	title	budget	revenue	Net Profit/Loss	year
11159	Chaos	20000000.0	10289.0	0.000514	2005.0
19027	5 Days of War	20000000.0	17479.0	0.000874	2011.0
21034	Special Forces	10000000.0	10759.0	0.001076	2011.0
25732	Foodfight!	65000000.0	73706.0	0.001134	2012.0
38388	Term Life	16500000.0	21256.0	0.001288	2016.0
19505	Laurence Anyways	9500000.0	12250.0	0.001289	2012.0
12038	The Good Night	15000000.0	20380.0	0.001359	2007.0
3966	Cherry 2000	10000000.0	14000.0	0.001400	1987.0
22097	Twice Born	13000000.0	18295.0	0.001407	2012.0
5651	All The Queen's Men	15000000.0	23000.0	0.001533	2001.0

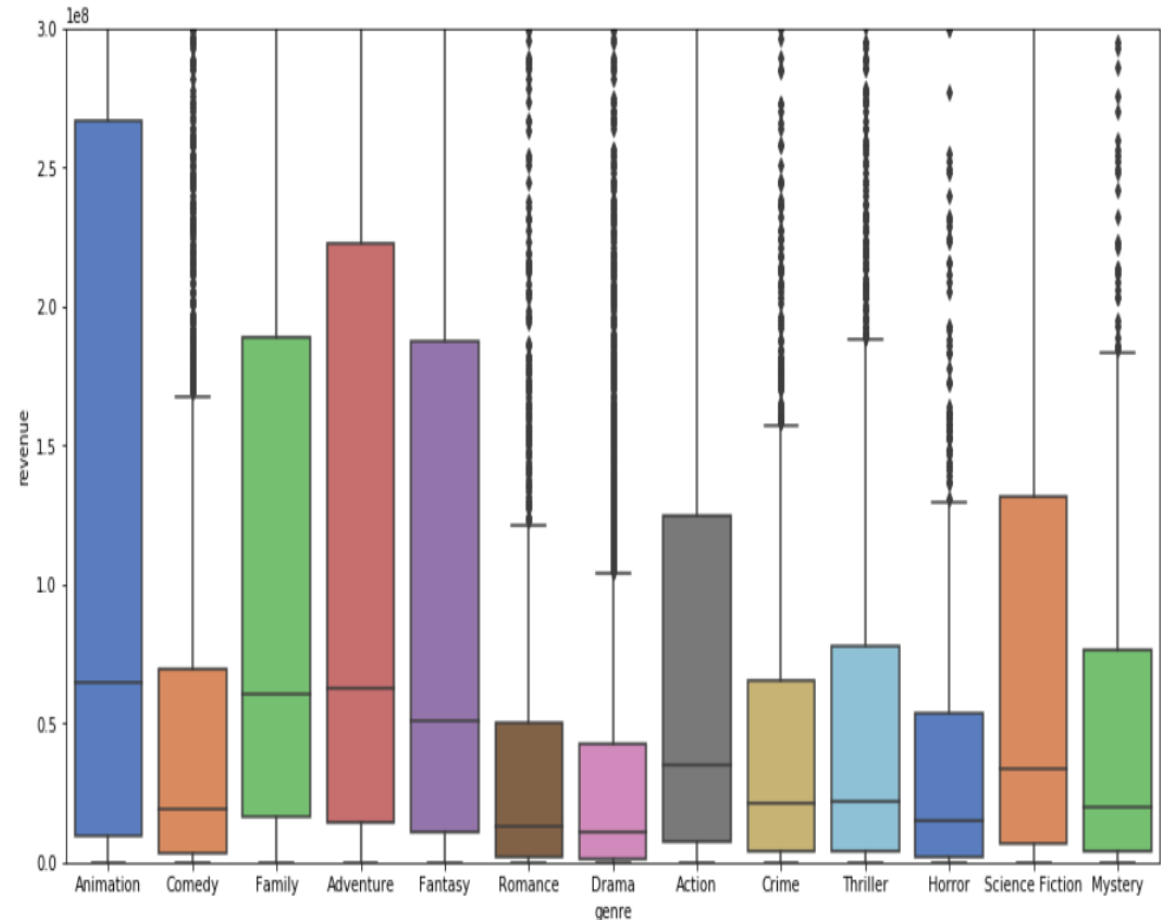
Most commonly occurring genres in movies

- TMDB defines 32 different genres for our set of 45,000 movies.
- Drama is the most commonly occurring genre, with almost half the movies identifying itself as a drama film.
- Comedy comes in at a distant second, with 25% of the movies having adequate humor doses.
- Other significant genres represented in the top 10 are Action, Horror, Crime, Mystery, Science Fiction, Animation, and Fantasy.



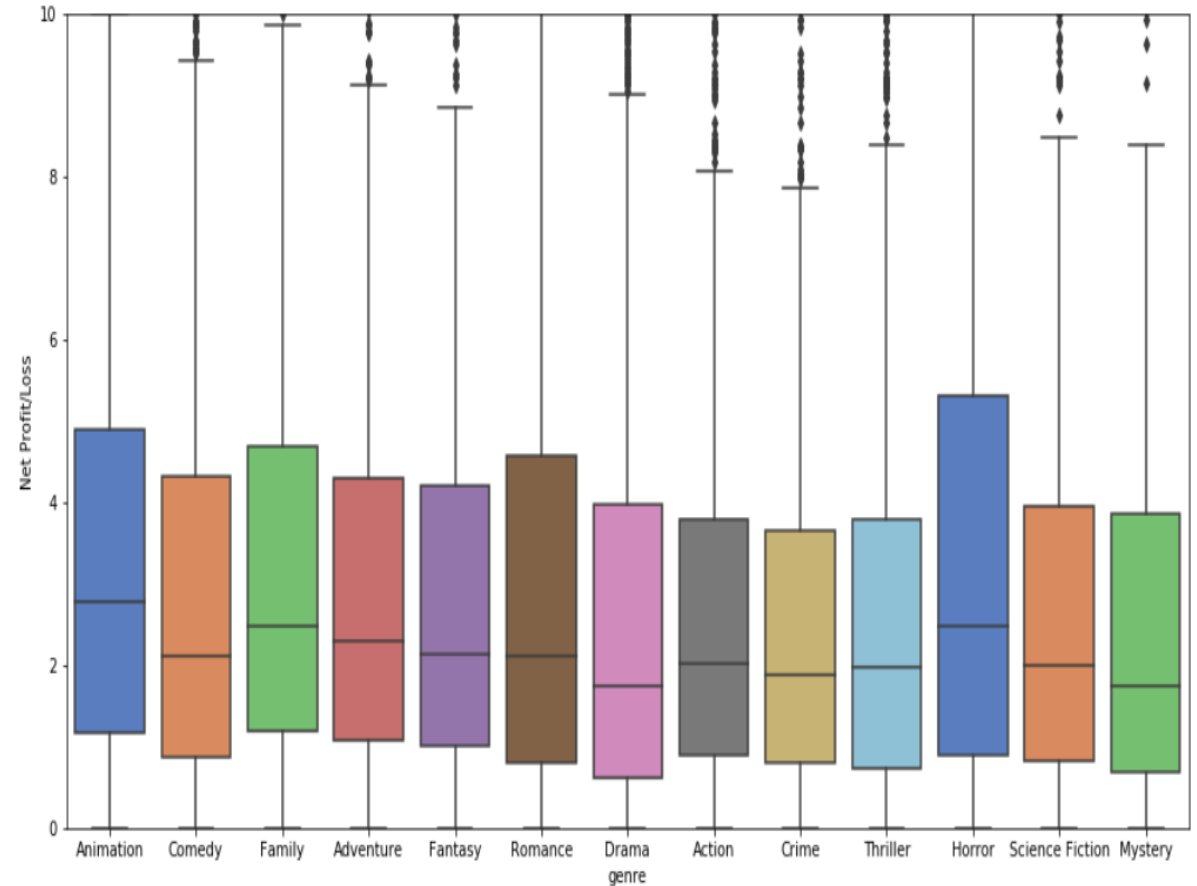
Are some genres more successful than others?

- Animation movies have the most extensive 25-75 range and the median revenue among all the genres plotted.
- Family and Adventure have the second and third highest median revenue, respectively.
- It is surreal to know that Romance and Drama have the lowest median revenue compared to other genres.



Are some genres more successful than others?

- Animation Movies tend to yield the highest returns on average. Horror Movies also tend to be a good yield.



Content-based recommenders:

- suggest similar items based on a particular item
- This system uses item metadata, such as genre, director, description, actors, etc. for movies, to make these recommendations.
- The general idea behind these recommender systems is that if a person likes a particular item, they will also enjoy an item that is similar to it.
- And to recommend that, it will make use of the user's past item metadata.

Collaborative Filtering recommenders:

- They try to predict the rating or preference that a user would give an item-based on past ratings and other users' preferences.
- Collaborative filters do not require item metadata like its content-based counterparts.

Simple recommender:

- The Simple Recommender offers generalized recommendations to every user based on movie popularity and genre.
- The basic idea behind this Recommender is that movies with more votes and more positive critic reception will have a higher probability of being liked by the average audience
- This model does not give personalized recommendations based on the user.
- In simple terms, this is like searching "Top 10 Good romantic movies of all time" in any search engine.
- I use the TMDB Ratings to come up with our Top Movies Chart. I will use IMDB's weighted rating formula to construct my chart.

Simple recommender continued..

- For a movie to feature in the charts, it must have more votes than at least 80% of the movies on the list.
- Therefore, to qualify for the chart, a movie has to have at least 50 votes on TMDB.
- We also see that the average rating for a movie on TMDB is 5.244 on a scale of 10, and 9151 Movies qualify to be on our chart.

Top 250 high rated movies Chart

	title	year	vote_count	vote_average	popularity	genres	weighted_Rating
10309	Dilwale Dulhania Le Jayenge	1995	661	9	34.457	[Comedy, Drama, Romance]	8.735928
15480	Inception	2010	14075	8	29.1081	[Action, Thriller, Science Fiction, Mystery, A...	7.990247
12481	The Dark Knight	2008	12269	8	123.167	[Drama, Action, Crime, Thriller]	7.988818
22879	Interstellar	2014	11187	8	32.2135	[Adventure, Drama, Science Fiction]	7.987741
2843	Fight Club	1999	9678	8	63.8696	[Drama]	7.985839
4863	The Lord of the Rings: The Fellowship of the Ring	2001	8892	8	32.0707	[Adventure, Fantasy, Action]	7.984595
292	Pulp Fiction	1994	8670	8	140.95	[Thriller, Crime]	7.984202
314	The Shawshank Redemption	1994	8358	8	51.6454	[Drama, Crime]	7.983616
7000	The Lord of the Rings: The Return of the King	2003	8226	8	29.3244	[Adventure, Fantasy, Action]	7.983355
351	Forrest Gump	1994	8147	8	48.3072	[Comedy, Drama, Romance]	7.983194

Top 250 high rated movies Chart

- It is interesting to see three Christopher Nolan Films, Inception, The Dark Knight, and Interstellar, which occur at the top of our chart.
- The chart also indicates a strong bias of TMDB Users towards particular genres and directors.

Top 10 drama movies

	title	year	vote_count	vote_average	popularity	weighted_Rating
12481	The Dark Knight	2008	12269	8	123.167	7.924623
22879	Interstellar	2014	11187	8	32.2135	7.917574
2843	Fight Club	1999	9678	8	63.8696	7.905213
314	The Shawshank Redemption	1994	8358	8	51.6454	7.890901
351	Forrest Gump	1994	8147	8	48.3072	7.888202
834	The Godfather	1972	6024	8	41.1093	7.851163
24860	The Imitation Game	2014	5895	8	31.5959	7.848105
359	The Lion King	1994	5520	8	21.6058	7.838458
18465	The Intouchables	2011	5410	8	16.0869	7.835391
22841	The Grand Budapest Hotel	2014	4644	8	14.442	7.810313

Top 10 Mystery movies

	title	year	vote_count	vote_average	popularity	weighted_Rating
15480	Inception	2010	14075	8	29.1081	7.856221
46	Se7en	1995	5915	8	18.4574	7.682311
11354	The Prestige	2006	4510	8	16.9456	7.598743
4099	Memento	2000	4168	8	15.4508	7.571292
9430	Oldboy	2003	2000	8	10.6169	7.243008
877	Rear Window	1954	1531	8	17.9113	7.092712
896	Citizen Kane	1941	1244	8	15.8119	6.967234
876	Vertigo	1958	1162	8	18.2082	6.924746
14825	Shutter Island	2010	6559	7	15.8136	6.822468
23675	Gone Girl	2014	6023	7	154.801	6.808585

Limitations of simple recommender

- It gives the same recommendation to everyone, regardless of the user's taste.
- If a person who loves romantic movies (and hates action) were to look at our Top 10 Chart, she/they/he wouldn't probably like most of the movies.
- If she/they/he were to go one step further and look at our charts by genre, she/they/he wouldn't still be getting the best recommendations.

Content based recommender

- Recommender based on the following metadata: the 3 top actors, the director, related genres, and the movie plot keywords.
- In this recommender system, the movie's content (cast, crew, keyword, etc.) is used to find its similarity with other movies.
- Then the movies that are most likely to be similar are recommended.
- As Required data was present in the form of "stringified" lists, I converted it into a safe and usable structure.
- I also converted the names and keyword instances into lowercase and stripped all the spaces between them.
- This is done so that our vectorizer doesn't count the Jennifer of "Jennifer Garner" and "Jennifer Aniston" as the same.

Content based recommender

- I have Mentioned Director 3 times to give it more weight relative to the entire cast.
- I created a "metadata soup," a string containing all the metadata that we want to feed to our vectorizer (namely actors, director, and keywords).
- To use textual data for predictive modeling, the text must be parsed to remove certain words – this process is called tokenization.
- These words need to be encoded as integers, or floating-point values, for inputs in machine learning algorithms. This process is called feature extraction (or vectorization).
- I have used **CountVectorizer()** for this process.

Content based recommender

- Using cosine similarity, I am calculating a numeric quantity that denotes the similarity between two movies.
- I used the cosine similarity score since it is independent of magnitude and is relatively easy and fast to calculate.

Algorithm for Content based recommender

- Get the index of the movie given its title.
- Get the list of cosine similarity scores for that particular movie with all movies.
- Convert it into a list of tuples where the first element is its position, and the second is the similarity score.
- Sort the list mentioned above of tuples based on the similarity scores, the second element.
- Get the top 10 elements of this list. Ignore the first element as it refers to self (the movie most similar to a particular movie is the movie itself).
- Return the titles corresponding to the indices of the top elements

Content based recommender-output

- The function **get_recommendations()** will take the movie title as the first input argument and the **cosine_sim** matrix as your second input argument.
- For this, we need a reverse mapping of movie titles and Data Frame indices.
- we need a mechanism to identify the index of a movie in our metadata Data Frame, given its title.
- The recommendations seem to have recognized other Christopher Nolan movies (due to the high weightage given to the director) and put them as top recommendations.

```
get_recommendations('The Dark Knight', cosine_sim)
```

```
8031    The Dark Knight Rises
6218           Batman Begins
6623           The Prestige
2085           Following
4145           Insomnia
7648           Inception
3381           Memento
8613           Interstellar
6645           Harsh Times
6902           Hitman
Name: title, dtype: object
```

Content based recommender-using vote count

- I took the top 25 movies based on similarity scores and calculated the 60th percentile movie's vote.
- Using this value of, we will calculate the weighted rating of each movie.

	title	vote_count	vote_average	year	weighted_rating
7648	Inception	14075	8	2010	7.990247
8613	Interstellar	11187	8	2014	7.987741
6623	The Prestige	4510	8	2006	7.969791
3381	Memento	4168	8	2000	7.967341
8031	The Dark Knight Rises	9263	7	2012	6.990577
6218	Batman Begins	7511	7	2005	6.988394
2839	American Psycho	2128	7	2000	6.959708
4145	Insomnia	1181	6	2002	5.969330
7912	Takers	399	6	2010	5.915913
6902	Hitman	982	5	2007	5.011865

Collaborative filter-based recommender

- To better interpret the data, we first pivot the data frame to have userId as rows and movieId as columns, filling the null values with 0.0.
- I have used the scipy library in Python to implement algorithms like Singular Value Decomposition (SVD) to give great recommendations and build a function that uses factorized matrices to recommend movies to a user user_id.
- This function evaluates all the movies and their ratings, rated so far by the user(impling the user watched this movie) and recommends other movies in the dataset by predicting its rating.

	title	vote_count	vote_average	year	weighted_rating
7648	Inception	14075	8	2010	7.990247
8613	Interstellar	11187	8	2014	7.987741
6623	The Prestige	4510	8	2006	7.969791
3381	Memento	4168	8	2000	7.967341
8031	The Dark Knight Rises	9263	7	2012	6.990577
6218	Batman Begins	7511	7	2005	6.988394
2839	American Psycho	2128	7	2000	6.959708
4145	Insomnia	1181	6	2002	5.969330
7912	Takers	399	6	2010	5.915913
6902	Hitman	982	5	2007	5.011865

Collaborative filter-based recommender

User 44 has already rated 25 movies.

Recommending the highest 10 predicted ratings movies not already rated.

User 44 has already rated 25 movies.

Below are 10 movies which user 44 has already rated.

```
already Rated.dropna().head(10)
```

	userID	MovieID	Rating	timestamp	Title	Genre
19	44	780	5.0	858707138	The Passion of Joan of Arc	[drama, history]
7	44	62	5.0	858707138	2001: A Space Odyssey	[sciencefiction, mystery, adventure]
23	44	805	4.0	858707310	Rosemary's Baby	[horror, drama, mystery]
9	44	104	4.0	858707248	Run Lola Run	[action, drama, thriller]
10	44	135	4.0	858707310	Dont Look Back	[documentary, music]
14	44	628	3.0	858707310	Interview with the Vampire	[horror, romance]
22	44	802	3.0	858707310	Lolita	[drama, romance]
21	44	788	3.0	858707248	Mrs. Doubtfire	[comedy, drama, family]
20	44	786	3.0	858707194	Almost Famous	[drama, music]
16	44	648	3.0	858707138	Beauty and the Beast	[drama, fantasy, romance]

Collaborative filter-based recommender

Recommending the highest 10 predicted ratings movies not already rated by user 44.

predictions

MovieID		Title	Genre
4168	608	Men in Black II	[action, adventure, comedy, sciencefiction]
2172	1073	Arlington Road	[drama, thriller, mystery]
1025	832	M	[drama, action, thriller, crime]
3211	708	The Living Daylights	[action, adventure, thriller]
1103	653	Nosferatu	[fantasy, horror]
5024	79	Hero	[drama, adventure, action, history]
923	762	Monty Python and the Holy Grail	[adventure, comedy, fantasy]
5437	673	Harry Potter and the Prisoner of Azkaban	[adventure, fantasy, family]
6287	647	Final Fantasy VII: Advent Children	[action, adventure, animation, fantasy]
6590	86	The Elementary Particles	[drama, romance]

Conclusion

- **Simple Recommender:** This system used overall TMDB Vote Count and Vote Averages to build Top Movies Charts, in general, and for a specific genre. The IMDB Weighted Rating System was used to calculate ratings on which the sorting was finally performed. This Recommender does not account for any user preference.
- **Content-Based Recommender:** I built a that took movie metadata such as cast, crew, genre, and keywords to make predictions. We also devised a simple filter to provide further preference to movies with more votes and higher ratings. However, It is only capable of suggesting movies that are close to a particular movie. That is, it is not capable of capturing tastes and providing recommendations across genres.
- **Collaborative Filtering:** We used the powerful Surprise Library to build a collaborative filter based on single value decomposition. One con about this Recommender is that it doesn't necessarily succeed in automatically matching content to one's preferences. These collaborative filtering systems require a substantial number of users to rate a new item before that item can be recommended.

Recommendation

- We can use content-based recommender to suggest movies based on a particular movie user watched.
- Collaborative filter-based recommender would be advisable to provide recommendations to user not only based on their watched movie list but also their likings towards it.
- With a hybrid model that brings together ideas from content and collaborative filtering to build an engine, we may be able to give better movie suggestions to a particular user in the future.



Thank You