

Data Exploration and Storyboard on the Dataset - IMDB 5000 Dataset

Group 16 - Akshat Tyagi (at3761), Harshit Srivastava (hs3500), Hamza Mirza (online, hm1800)

Data Set

The dataset contains information about movies that are listed on IMDB, aggregated at the movie level. So, it contains information in the following manner:

Director - Name of the director (categorical)

Actor - Names of 3 actors (categorical)

Facebook likes (0-1M) - director, actor and total cast (quantitative)

Genre - movie genre (categorical)

No. of user reviews (0-5000) - no. of reviews submitted by users (quantitative)

No. of critic reviews (0-800) - no. of reviews submitted by critics (quantitative)

IMDB rating (0-10) - IMDB rating considering both user reviews and critics' reviews (quantitative, ordinal)

Duration (0-500) - Duration of the movie in mins (quantitative, ordinal)

Title_year - Year when the movie was released. (quantitative, temporal)

Language - Language in which movie was released. (categorical)

Budget - Amount of money spent on making a movie. (quantitative)

Gross - Amount earned by the money during its run. (quantitative)

Movie_title - The name of the movie. (Categorical)

Plot_keywords - Some words used to describe a movie (Categorical).

Dataset Link - (<https://www.kaggle.com/carolzhangdc/imdb-5000-movie-dataset>)

Analytical Questions and Proxy Tasks

We will explore the IMDB movie reviews dataset to answer some pertinent questions regarding how movies are rated based on a plethora of factors. We would like to answer the following questions -

Q1. Who is the best-rated director?

Proxy Task: Analyzing distribution of movie ratings for different directors

Proxy Value: best-rated -> high IMDB rating

Q2. Which are the best-rated movies?

Proxy Task: Visualizing features of best-rated movies

Proxy Value: best-rated -> high IMDB rating

Q3. Does the movie rating change with time?

Proxy Task: Does the IMDB rating change with title_year?

Proxy Value: movie rating -> IMDB Rating, time -> title_year

Q4. Are movies in some languages better rated than movies in others?

Proxy Task: Visualizing rating distribution for movies of different languages

Proxy Value: better rated -> high IMDB rating, language -> language

Q5. A movie is a hit if it earns more money than it costs. What are some of the biggest hits? Do some directors have more hits than others?

Proxy Task: Do movie_titles from some directors have a higher difference between budget and gross than others?

Proxy Value: movie -> movie_title, director -> director, earns -> gross, costs -> budget.

Q6. Does the plot keyword of movies change with time?

Proxy Task: Do aggregated plot_titles for the movie_title change with title_year?

Proxy Value: plot keyword -> plot_titles, movie -> movie_title, time -> title_year

Q7. Does the movie genre have any effect on the movie's rating?

Proxy Task: What is the average movie rating for a given genre? Do any specific genre has a higher average rating compared to other genres?

Proxy Value: higher average rating -> avg(IMDB Rating), genre -> genres

Q8. Does the budget of the movie relate to the IMDB rating?

Proxy Task: What is the average IMDB rating for a given budget range?

Proxy Value: Average IMDB rating -> avg(IMDB rating), budget range -> budget

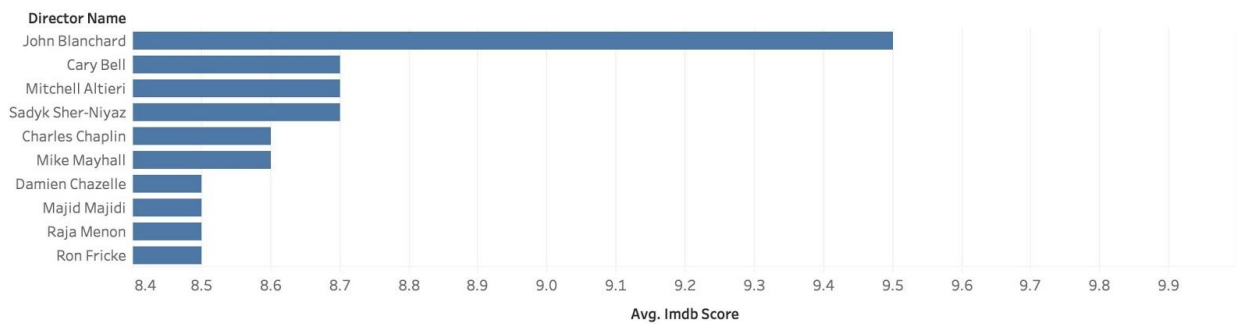
In order to answer these questions, we would like to create interactive visualizations where we can change the visualized data based on the attributes. The website would consist of a dashboard with multiple visualizations answering the above questions. We would like our website to be dynamic such that it can process the data for the visualizations on the backend in less time. For this, we would query our data with cross filter, which would be used for the visualizations. With the interactive nature of the visualizations, not only will we answer some important questions about the dataset, but we would also like the user to interact with the visualizations and form their own questions.

Story Design

Data Analysis

Q1. Who is the best-rated director?

Best Rated Director

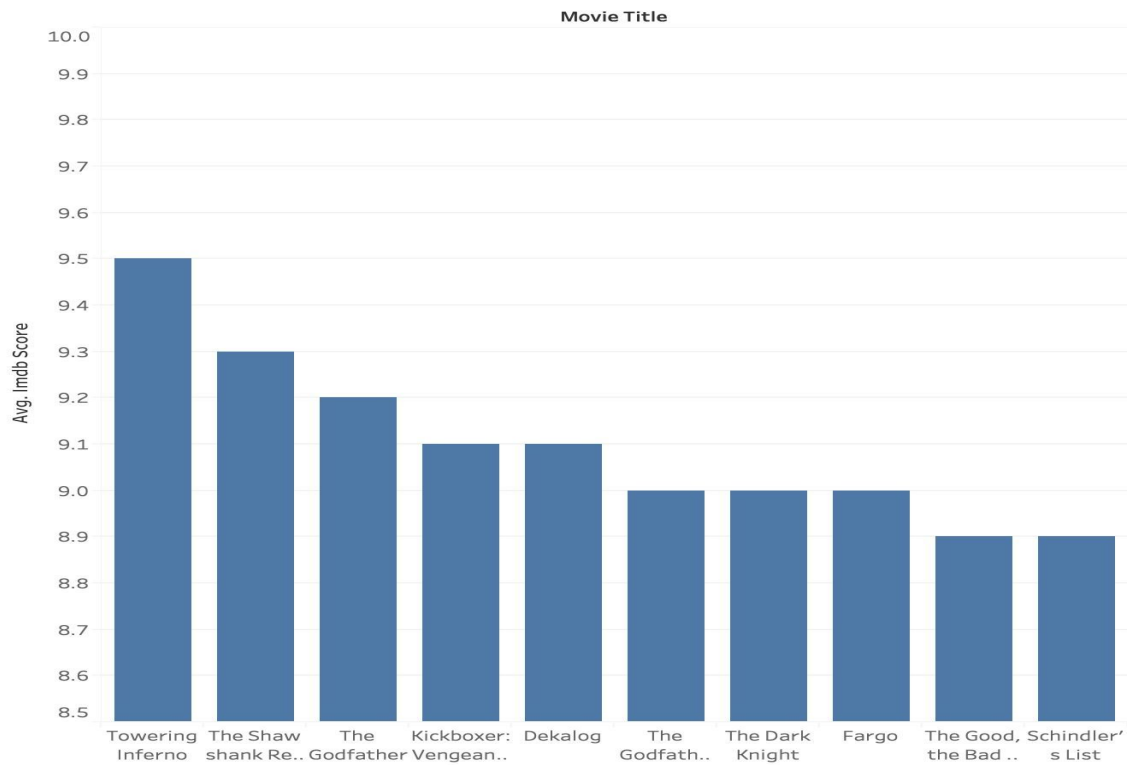


Average of Imdb Score for each Director Name. The view is filtered on Director Name, which has multiple members selected.

It was observed that John Blanchard , Cary Bell, and Mitchell Altieri are the best directors because their movies have the highest ratings.

Q2. Which are the best-rated movies?

Best Rated Movies

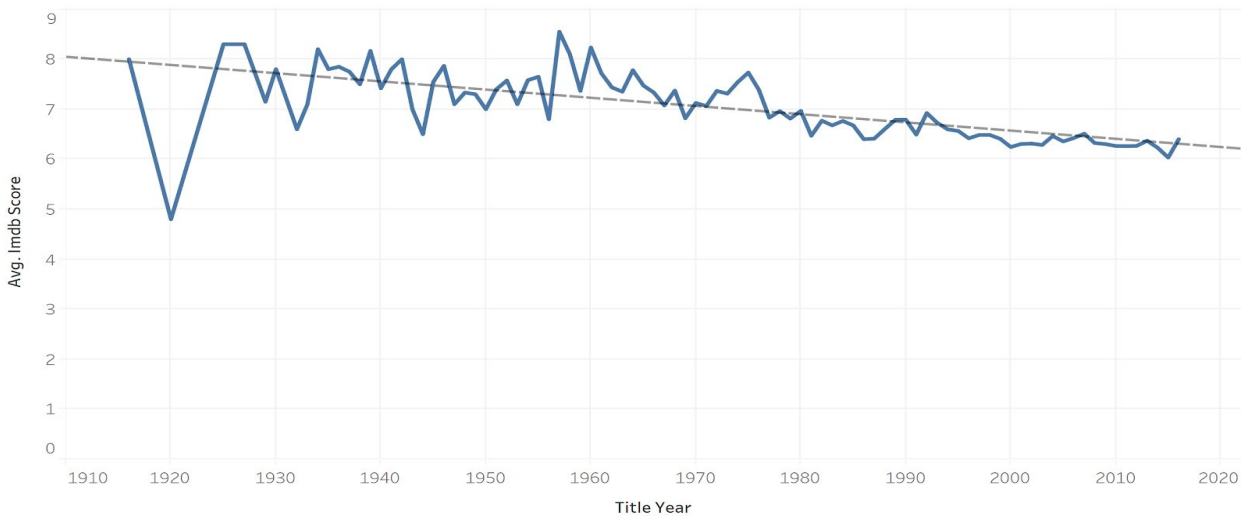


Average of Imdb Score for each Movie Title. The view is filtered on Movie Title, which has multiple members selected.

Towering Inferno, The Shawshank Redemption, and The Godfather are the highest rated movies.

Q3. Does the movie rating change with time?

Movie Ratings vs time

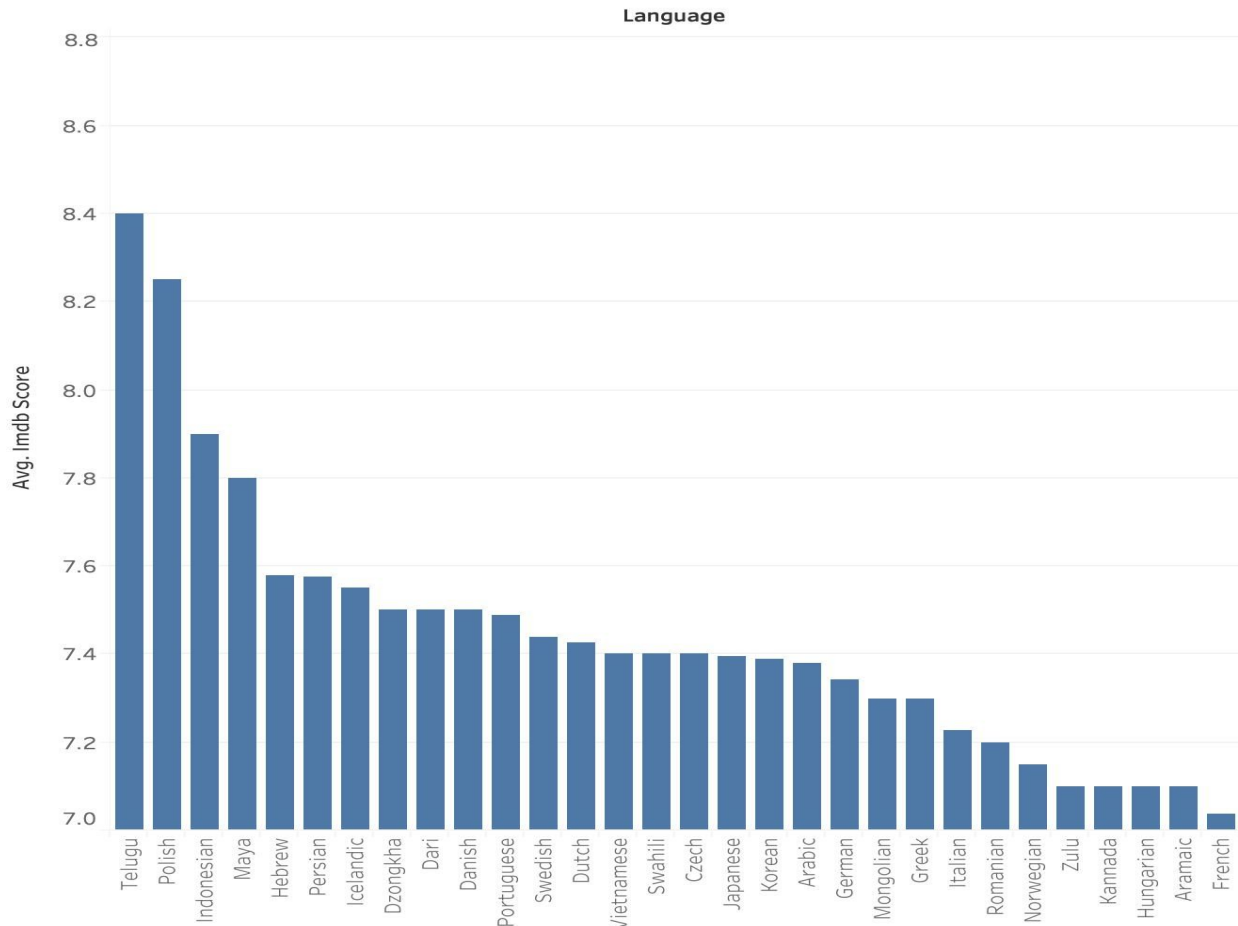


The trend of average of Imdb Score for Title Year.

Yes, it was observed that time time, the average IMDB rating of a movie decreased, as shown above.

Q4. Are movies in some languages better rated than movies in others?

Movie Language vs Rating

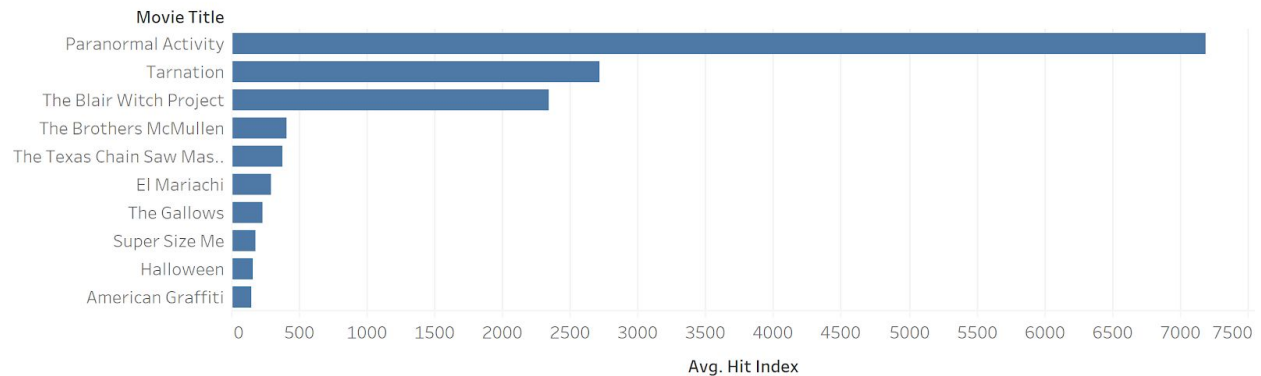


Average of Imdb Score for each Language. The view is filtered on Language, which keeps 30 of 48 members.

The above graph shows top 30 highly-rated languages. As shown above, it's clear that while movies in Telugu, Polish, and Indonesian have high ratings, movies in Chinese, Tamil, and Bosnian have low ratings. Thus, movies in some languages are rated better than others.

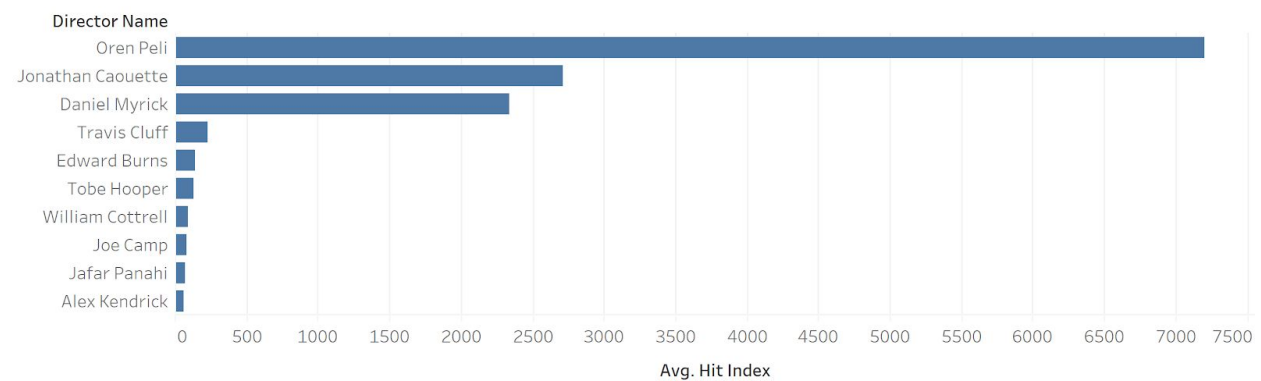
Q5. A movie is a hit if it earns more money than it costs. What are some of the biggest hits? Do some directors have more hits than others?

Biggest Hits



Average of Hit Index for each Movie Title. The view is filtered on Movie Title, which keeps 10 of 4,916 members.

Most Hit Directors



Average of Hit Index for each Director Name. The view is filtered on Director Name, which keeps 10 of 2,399 members.

Based on the analysis, Paranormal Activity is the biggest hit, followed by Tarnation and The Blair Witch Project. Oren Peli, Jonathan Caouette, and Daniel Myrick are the highest rated directors.

Q6. Does the plot keyword of movies change with time?

There was no change observed.

We have decided to omit this question as its not pertinent to our storytelling and we were rightly advised by the TAs to not consider this question.

Wordcloud

Plot Keywords. Color shows details about Plot Keywords. Size shows count of Plot Keywords.

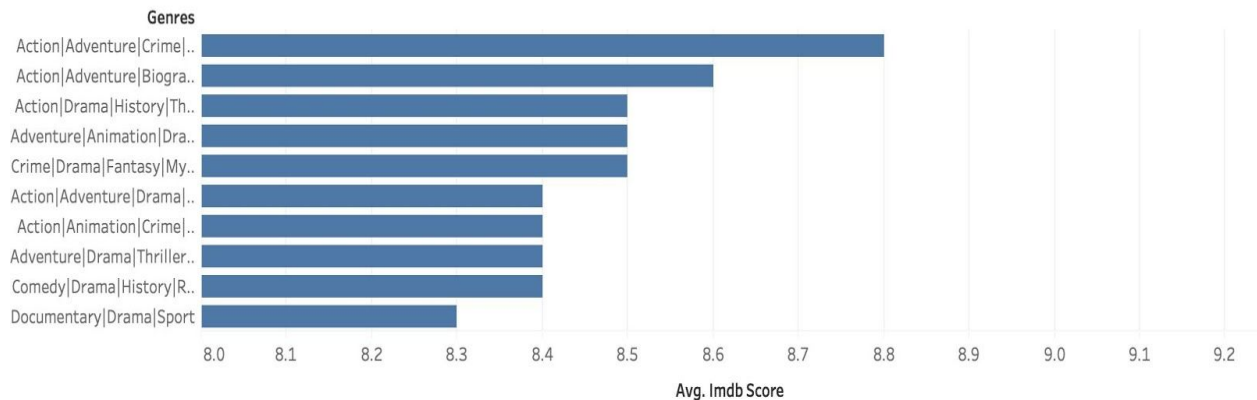
automobile association of america
 friend relationship|southern california|voice over narration directed by star|friendship|reference to isaac newton|skinny dipping|written and directed by cast member. based on comic
 best friend|black and white|father daughter relationship|female protagonist|mother daughter relationship true crime cosmetics|female protagonist|native american|native a
 brother relationship|high school friends|imperative in title|skateboard|brother sister relationship|family relationships|middle east|mother son relationship|twin brother and sister
 ages in end credits male objectification|manipulative behavior|modern day adaptation|narcissistic woman|promiscuous woman|hairsaw|chainsaw murder|human monster|music
 s|reference to vincent van gogh|talking about suicide female frontal nudity|female rear nudity|title directed by female|virgin|written and directed by cast member. aunt niece re
 frontal nudity|nude with glasses killed in police car|mercenary|police officer shot in the leg|police officer shot through the heart|whistleblower argument|reference to van ill
 stage play|freedom|historical fiction|virgin reference to ben affleck|reference to brad pitt|reference to george clooney|reference to jack nicholson|website|clone|cloning|construct
 on in a car trunk|school janitor full frontal male nudity|hotel room|multiple characters voiced by same person|one night stand|sex in hotel room 3 dimensional|concert footage|engl
 |subway tunnel|chrysler building manhattan new york city|delusion|manhattan new york city|new york city|new york city|title directed by female inspired by radio program|inspired by true event
 d by female older woman younger woman relationship|sex with student|teacher|teacher-student relationship|title directed by female animate object|imaginary friend|reference to
 ve brother in law sister in law relationship|expectant father|hi|chronic personality disorder|middle aged woman|off screen rape reference to snow white|sleeping bag|stepmotl
 personality disorder|unhappy marriage|widow alternate timeline|same actor playing two characters simultaneously on screen|second part|year 2015|younger version of character
 south africa british|character says now that you're on the case i hope we're going to have some gratuitous sex and violence|cruise missile|missile|nuclear warhead
 york city|new york city|urination modern day adaptation|reference to william shakespeare|shakespeare in modern dress|shakespeare play|shakespeare's coriolanus abusive husbar
 rist|mother daughter relationship|murder police officer killed|police officer shot in the chest|police officer shot in the forehead|police officer shot in the head|police shootout co
 ew york city american actress playing british character|book publishing|calorie counting|employer-employee relationship|rabbit costume body swap|camera shot of feet|foot close
 |post apocalypse|protective father|zombie apocalypse ex convict|graduation|manhattan new york city|older brother is bad influence|washington heights manhattan new york city t
 ralzheimer's disease|bechdel test passed|linguistics professor|new york city|reference to angels in america the play 3d railroad|reference to douglas fairbanks|reference to jesse
 merican reservation martial arts|murdered before giving protagonist information|part computer animation|prequel|prequel to cult film cartoon on tv|reference to frankenstein|
 rt|singing in a car actor playing multiple roles|brownie the creature|closing credits sequence|family relationships|magical creature apostrophe in title|critically bashed|hit on t
 |life father and son playing father and son full frontal male nudity|man wearing a jock strap|man wearing a thong|rear male nudity|stripped prison uniform fish out of water|j
 ighter relationship|reference to superman|super villain bloopers during credits|breakfast cereal|department store ride|face painting|sponge bob square pants slippers atrocity|ce
 ger bitten off museum|museum of natural history manhattan new york city|night watchman|star died before release|tablet breaking the fourth wall|breaking the fourth wall by talk
 |reference to howard university execution by hanging|los angeles police department|miscarriage of justice|wrongful conviction abusive relationship|character name li
 ntal nudity|male nudity|actress playing herself|actress shares first name with character|fired from a job|male nudity|written by star french revolution|old testament|part of an unfin
 eality apartment building|character's point of view camera shot|fire station|subjective camera|television reporter dwegons
 ckedg soles|foot closeup four word title|grandparent grandchild relationship|reference to facebook|reference to the internet|singer female protagonist|high school|self destruct
 ||librarian|male objectification|pregnant woman in bathtub based on young adult novel|dystopia|genetic experimentation|strong female character|strong female lead dog hit by a

Q7. Does the movie genre have any effect on the movie's rating?

Movies with tags such as Action, Adventure, Crime etc have better ratings than ones with the tags - Documentary, Drama, Sport.

We have decided to omit this question because of multi-class genres present in the data and there's no good way to aggregate it, on top of it, it doesn't really help with the story-telling.

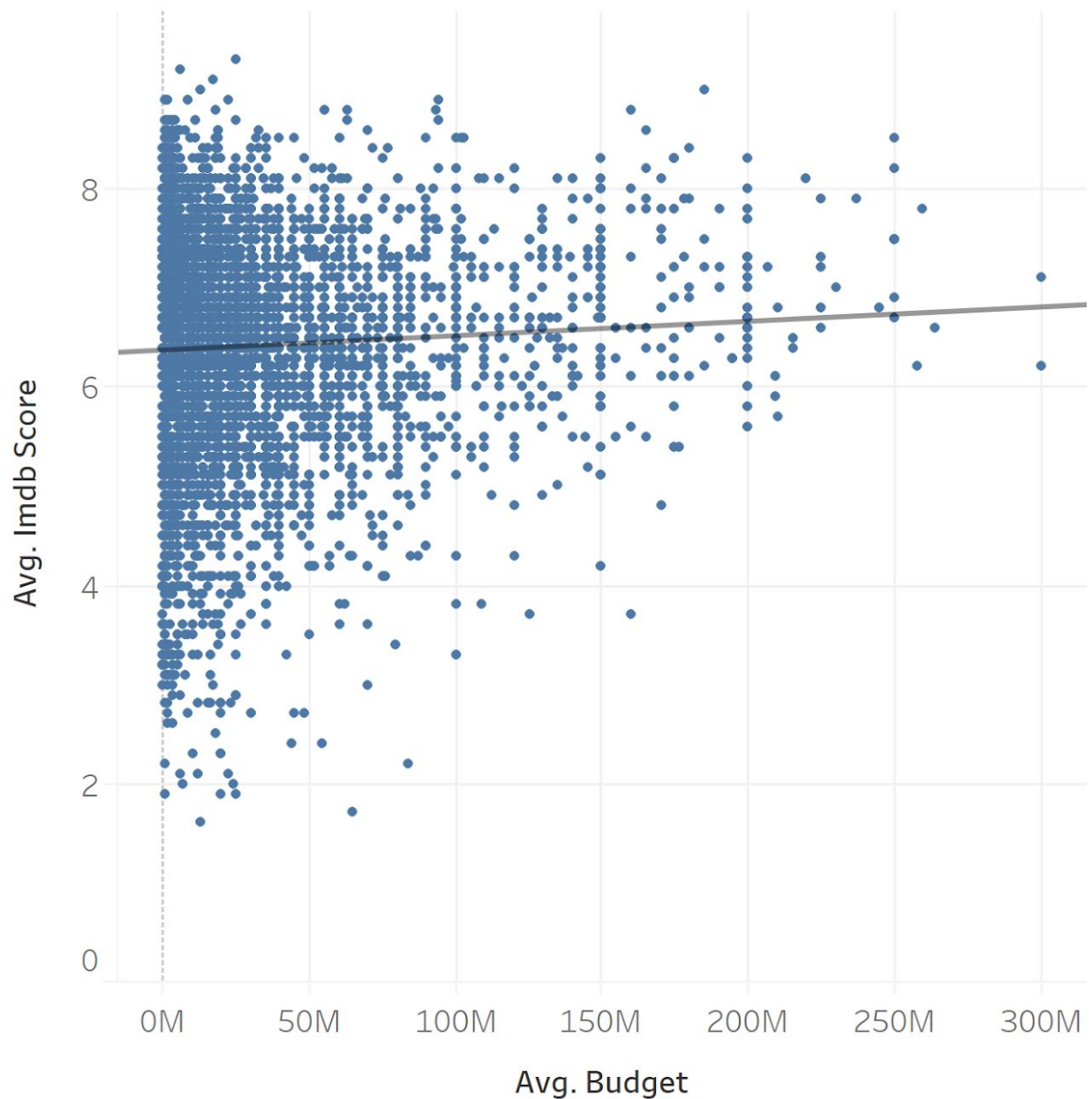
Sheet 9



Average of Imdb Score for each Genres. The view is filtered on Genres, which keeps 10 of 914 members.

Q8. Does the budget of the movie relate to the IMDB rating?

Movie Budget vs Rating



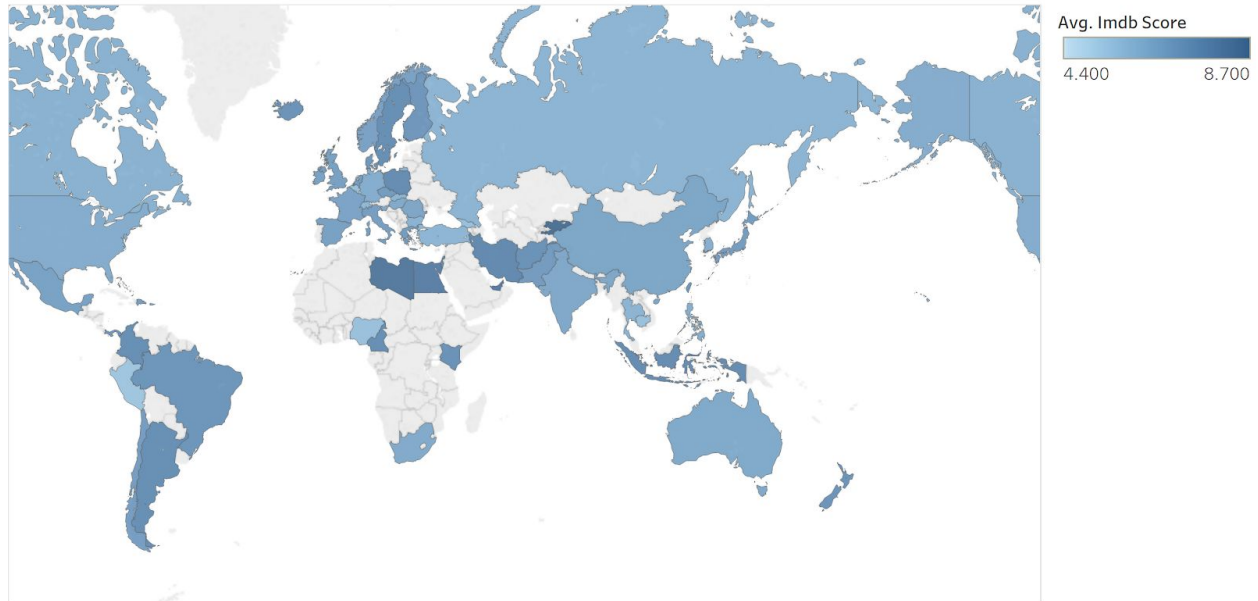
Average of Budget vs. average of Imdb Score. Details are shown for Movie Title. The view is filtered on average of Budget, which ranges from 218 to 323,253,073.

Yes, as shown above, the average imdb score and average movie budget have a positive correlation.

Additional Questions Asked after observing the dataset -

Q9. Does movie rating vary by country?

IMDB Ratings by country

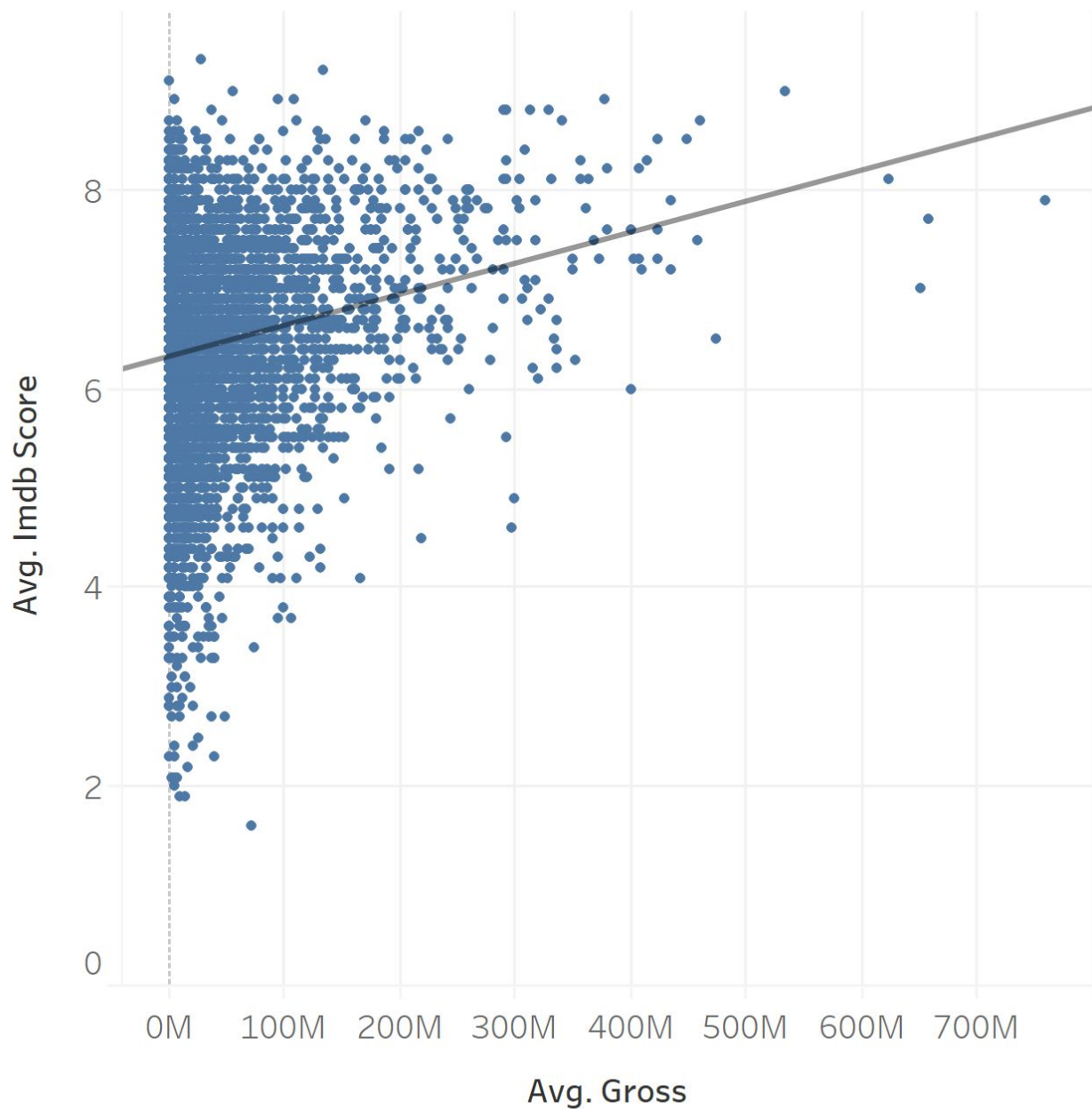


Map based on Longitude (generated) and Latitude (generated). Color shows average of Imdb Score. Details are shown for Country.

Yes. Countries like Iran, Egypt, and Libya have better rated movies than Russia, Nigeria etc.

Q11. Is the money a movie makes relating to its rating?

Average Grossing by IMDB Rating



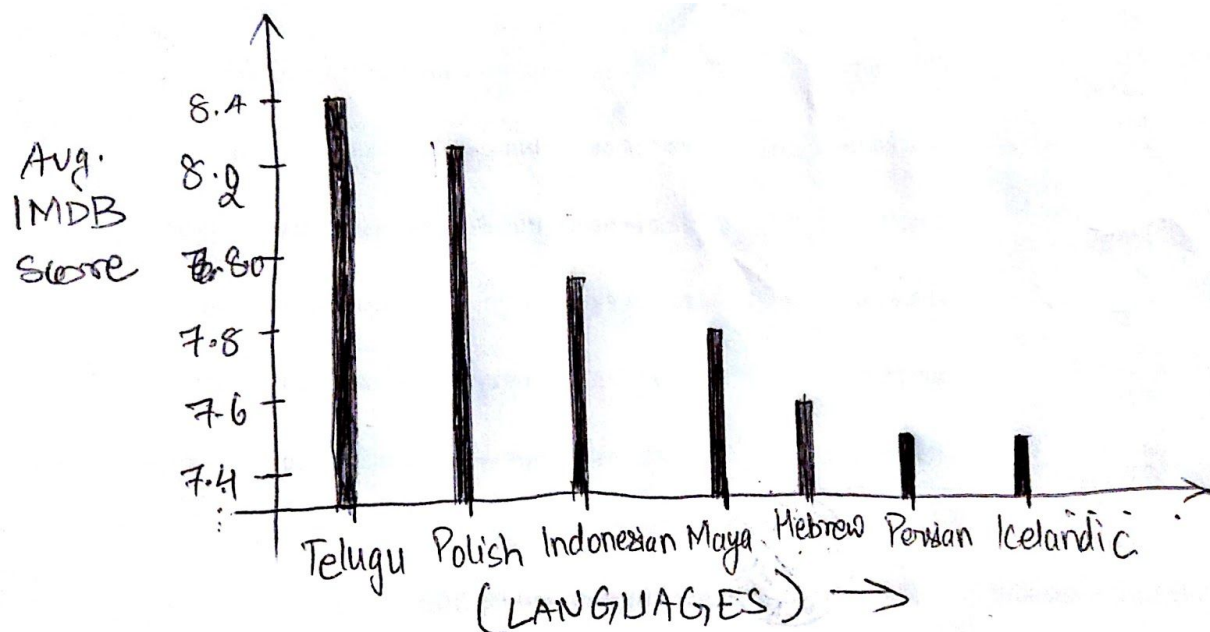
Average of Gross vs. average of Imdb Score. Details are shown for Movie Title.

Yes. Movies which have a good IMDB rating tend to earn more, as shown above.

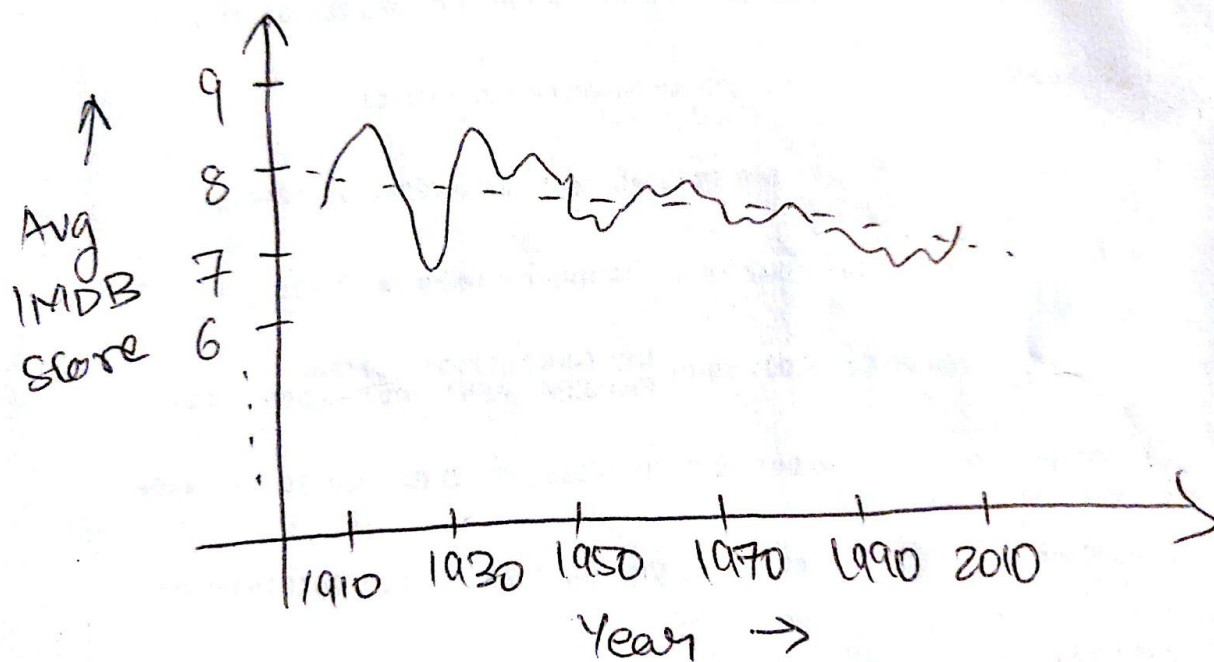
Storyboard

The IMDB dataset is fascinating because it is allowing us to observe some general trends about movies. A movie's performance depends on a plethora of factors, and we will observe a few of them today.

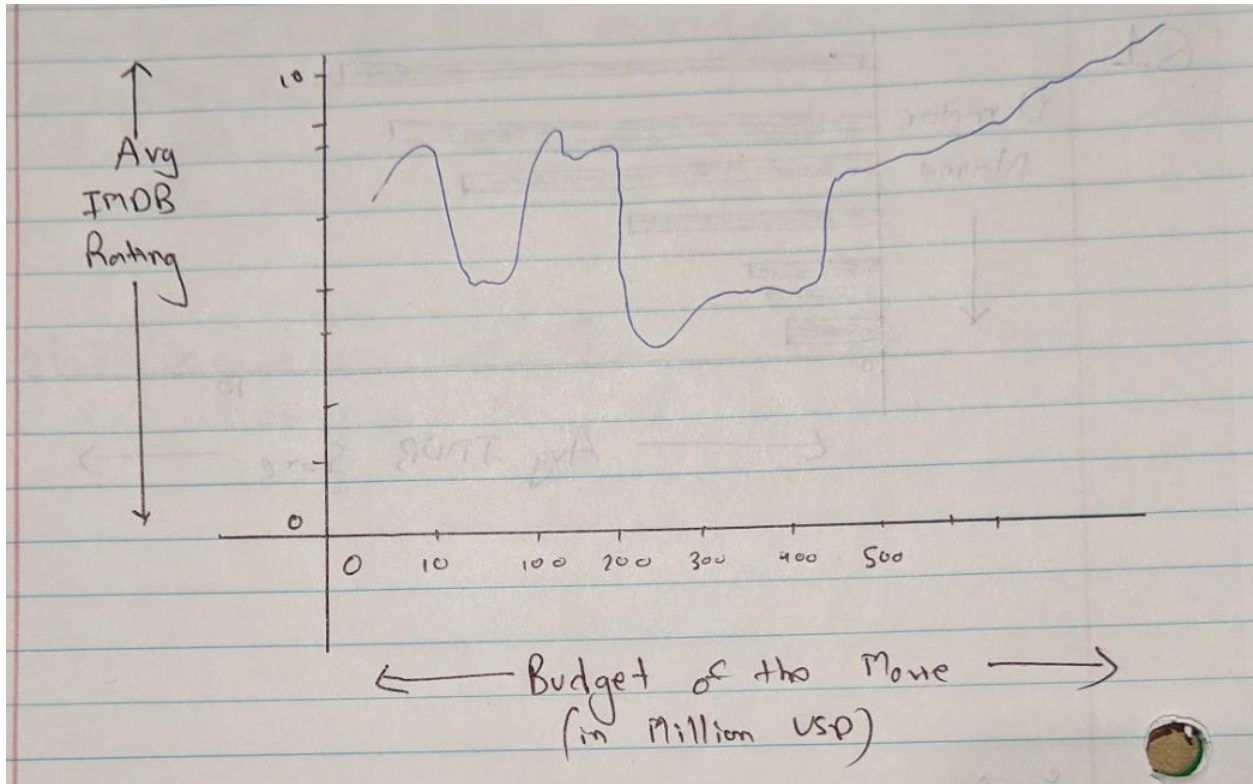
First, we observed that it's not just the content of the movie, but also its origins that affect its rating. As shown below, we can see that movies from some parts of the world and some languages do better than others.



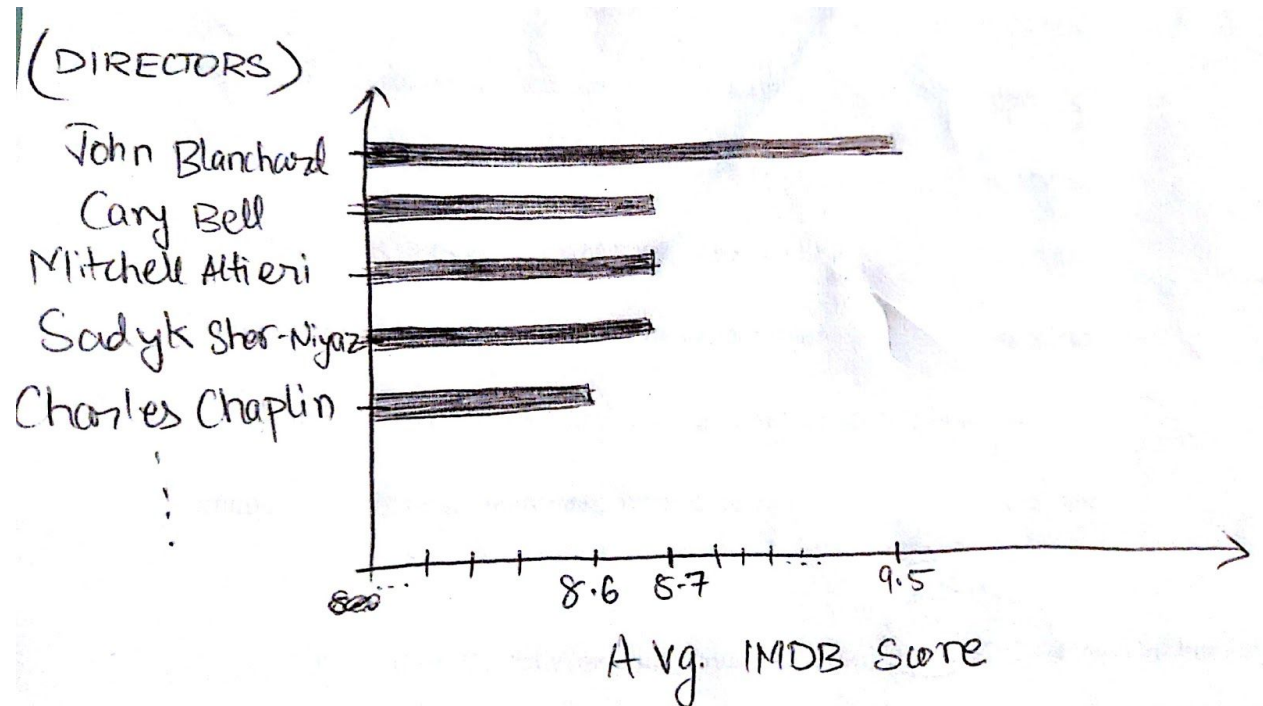
People mostly say that movies made in the olden days were better than the ones made today. We observed this hypothesis and concluded that the hypothesis holds true, based on the ratings of movies based on their release date, as shown below.



Now that we have established that a movie's origins plays a crucial role in its rating, we would also like to see if spending more money means that the movie is good. As shown below, the amount of money spent on a movie affects its rating.

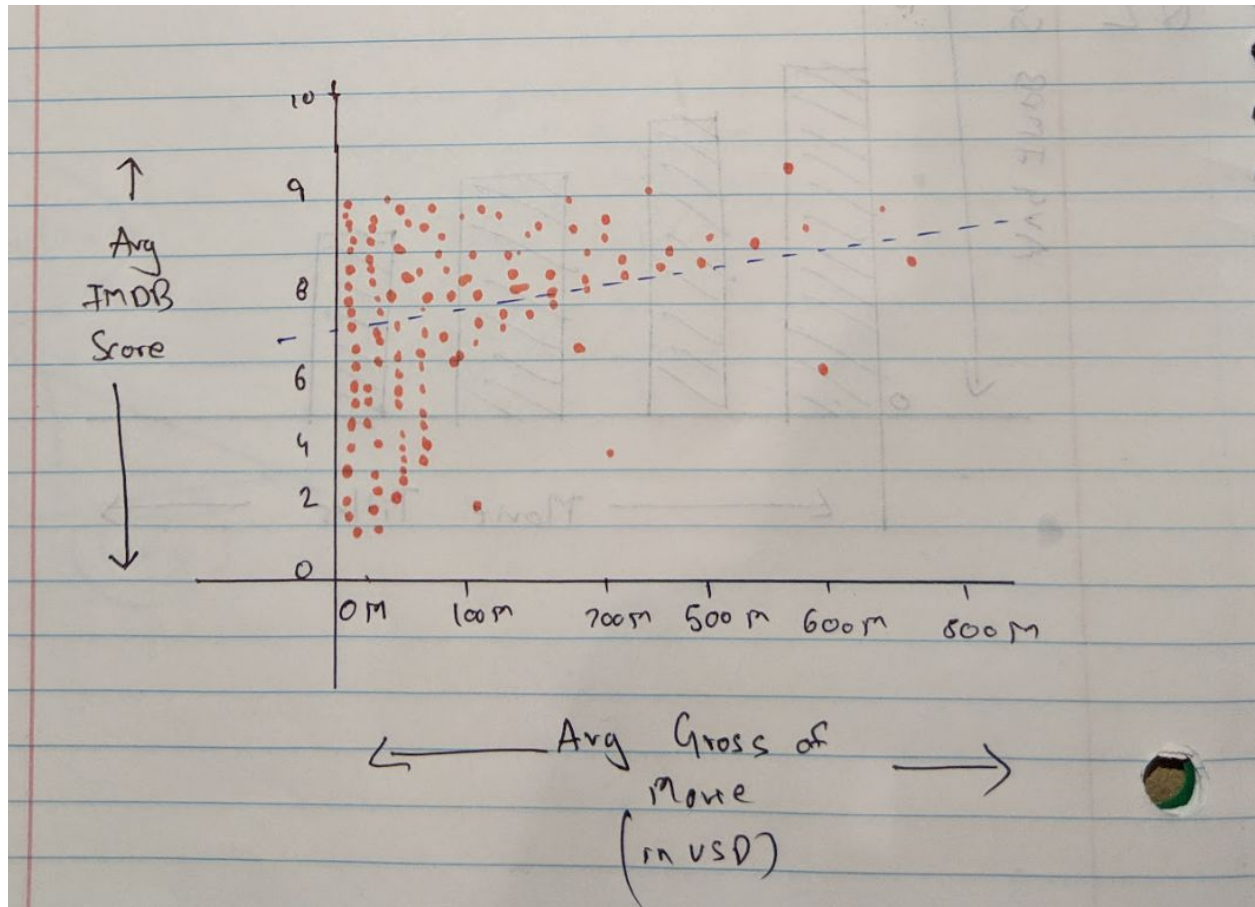


Since how well a movie is made depends a lot on directors, let's look at some of the best directors out there, based on how hit their movies were -

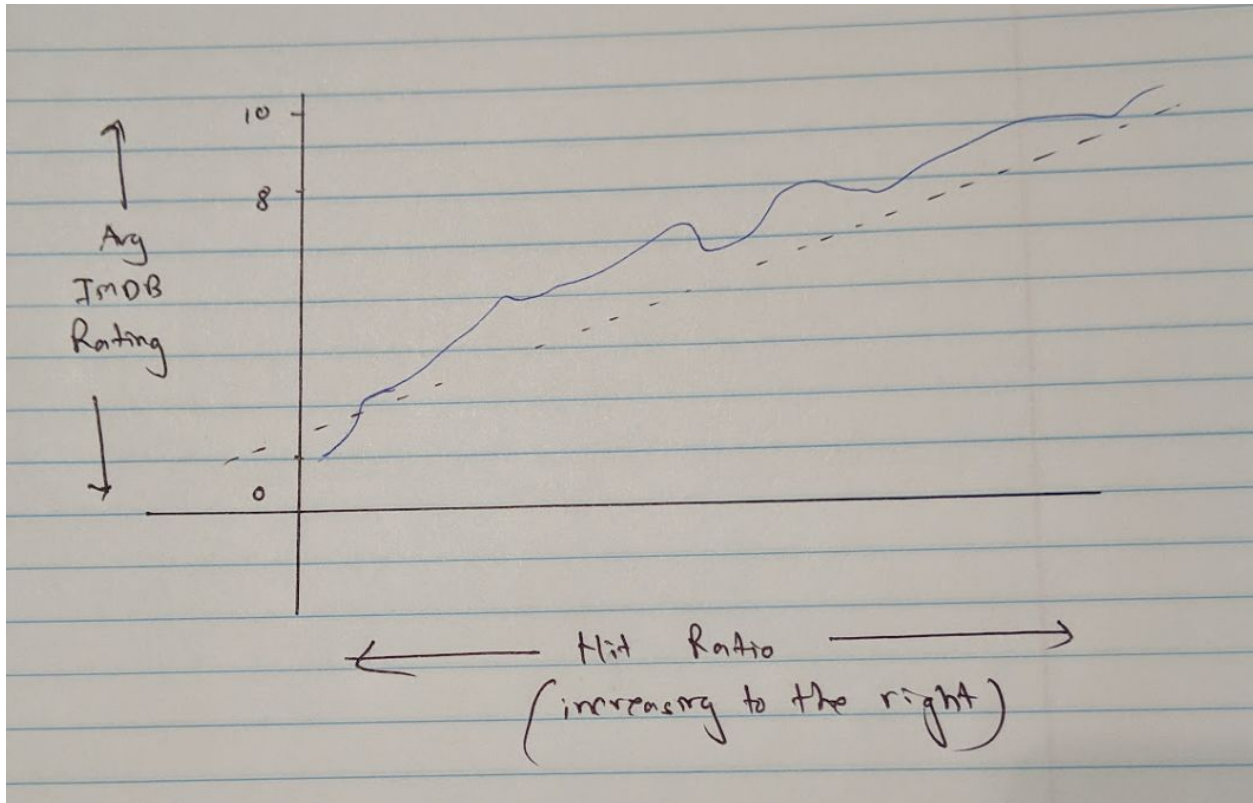


So far, we have observed indicators that could potentially predict the performance of a movie. Now let's look at some indicators which quantify how good or bad a movie is.

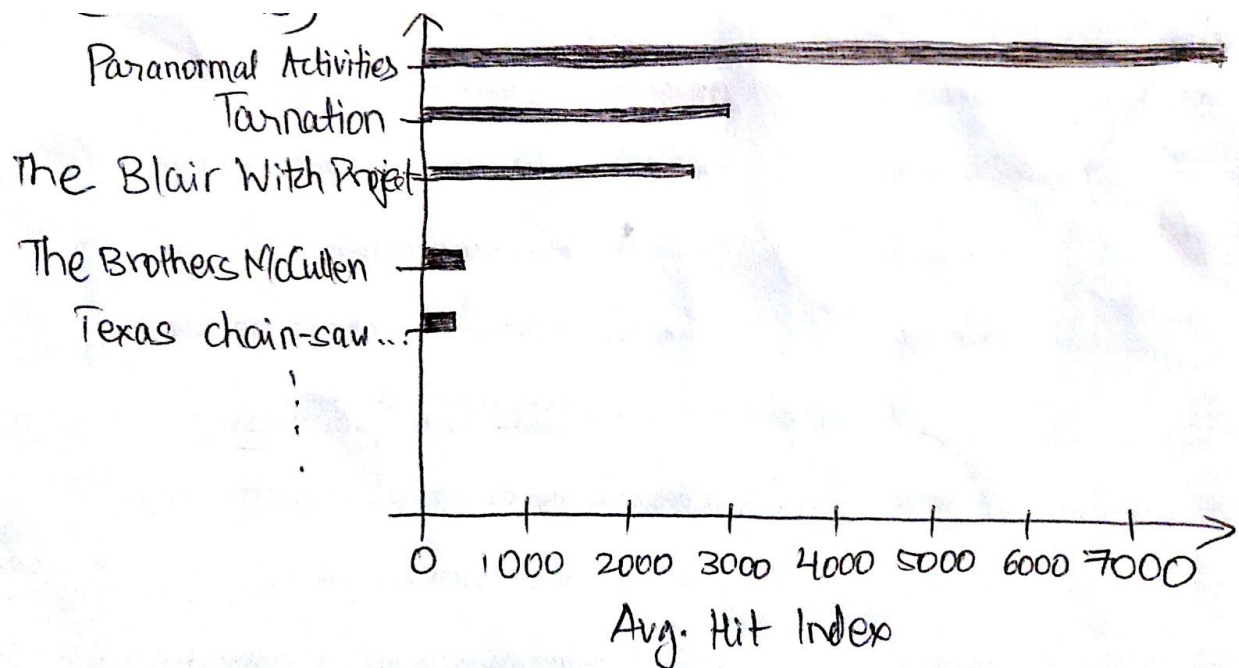
By observing the visualization below, it's clear that the money a movie makes is proportional to its IMDB rating. This means that there is a general trend where audience appreciates well-made movies.



Now let's see if movies that are hit, are hits because they're well made or not. We can observe from the graph below that well-rated movies turn out to be hits. Some, like Paranormal activity make a lot of money for a small budget.



To recapitulate, let's look at the highest rated movies of all times -



We can conclude that a movie's success is heavily dependant on its rating. But it's also dependant on a plethora of factors, including the language its made in, time when it was released, and the region it comes from.

Changelog

We answered two additional questions which we felt were germane to the analysis. Those questions are -

Q1. Does the average movie rating vary by country?

Q2. Does the amount of money a movie earn depend on its rating?

Phase 1 Update

By now, we have a working web page implemented. We have added the text from our stories and are currently working on adding the visualizations. We have added three new choropleth maps, with the possibility of adding a fourth, in addition to the original proposal. The new maps are as follows - Average movie budget for each country, average number of facebook likes per country, average number of votes per country. We are deliberating on adding a fourth choropleth map to visualize the average budget per country.

LIST OF IMPLEMENTED ITEMS -

1. Choropleth Map to visualize the average IMDB score grouped by the country of origin.
2. Choropleth Map to visualize the average movie budget grouped by the country of origin.
3. Choropleth Map to visualize the average number of facebook likes for the movie grouped by the country of origin.
4. Choropleth Map to visualize the average number of votes for the movie grouped by the country of origin.
5. Helper functions to load specific data, group them according to some feature, and scale them up.
6. Helper functions to create choropleth maps and assign specific properties to them.
7. Helper functions to create bar charts.

8. Bar chart for language vs imdb_score.

LIST OF ITEMS TO BE IMPLEMENTED

1. Helper functions to create scatter plots.
2. Creation of remaining bar charts and scatter plots using the helper functions for other visualizations in the project.
3. Making the visualizations interactive by adding tooltips, animations, filters etc.

CURRENT VERSION

<https://github.com/NYU-VIS-FALL2018/storytelling-group-16>