

Does a movie's rating depend of the choice of genre?

Karola Takács

02/01/2021

1. Executive Summary

I was interested in the connection of weighted average IMDB ratings and the genre of movies: is there a relationship indicating that some genres are more likely to be upvoted than others? The IMDB dataset is historically extensive however might be problematic since the same user base voted for old and for new films. Nonetheless the results are based on data from 1920 till 1995 on films which have more than 25 thousand votes to retain credibility on such subjective matter as taste. The analysis suggests that topics of Family, Horror or Fantasy tend to have lower ratings so screenwriters might consider choosing more 'likeable' themes or a director might decide for a book to adapt on screen which is not in these categories. I have also found that metacritic score is strongly correlated with IMDB ratings but has not much causal relationship with it.

2. Introduction

The primary interest of this analysis was whether there is a causal connection between the IMDB rating of a film and it's genre. I was simply interested if certain genres are more 'likable' by viewers or more successful from a director's point of view. My original (broad) research question was: "Is there a favorable genre to achieve higher movie rating?" And also "What other factors might influence the rating?" However, after the data cleaning part, I managed to narrow down the research question to: "Does the genre choice of a director yield to certain IMDB ratings for movies older than 1920? (more precisely: for movies with more than 25000 reviews). With this information at a director's hands it would be a direction what are the seemingly "more successful" genres.

3. Data

The main variables of the joined dataset to be analyzed are the movies' title, year of making, their genre along with ratings from two different sources: IMDB and Metacritic (Metacritic aggregates movie reviews from the leading critics) and also number of reviews by IMDB users and critics and finally the worldwide gross income in USD millions. The left-hand sided variable is rating and the right-hand side variable is genre.

To begin with I had two tables about movies both downloaded from the Kaggle webpage. They differed in record count, but after joining them based on the IMDB movie ID field, I had 43605 observations. The year variable ranges from 1874 to 1995.

During the cleaning process I excluded all the observations with missing values for the most important variables, IMDB rating and genre, but these accounted only for 2031 observations. Let me explain these two variables in more detail.

- IMDB rating: all registered members of IMDB can cast their votes. For one movie a user can vote several times but then the last value is going to be updated, so basically there is 1 user for 1 film at a single time (good data quality).

- Genre: this is a problematic variable since it not only is categorical but for a single movie there can be more genres attached. For simplification my assumption was: the firstly listed genre is the most relevant or primary genre. Even if this is true, it really makes sense to have more categories, because there are few main categories, like action or drama and it is easy to put most of the movies into one of these which results in that these categories will be overrepresented without other genres to balance out.

Descriptive statistics show that IMDB rating values are close to normal distribution, bit skewed to the left, signaling that the mean is well above 5, it is around 7. Similar with metacritic score, which has a better looking bell-shape, the distribution of values are wider, values are more spread on the scale of 100. IMDB Votes would make sense to transform for the analysis part as it shows right skewness, but since it is already incorporated into the IMDB rating variable it would be a bad conditioning variable in the visual analysis, so I decided not to use it.

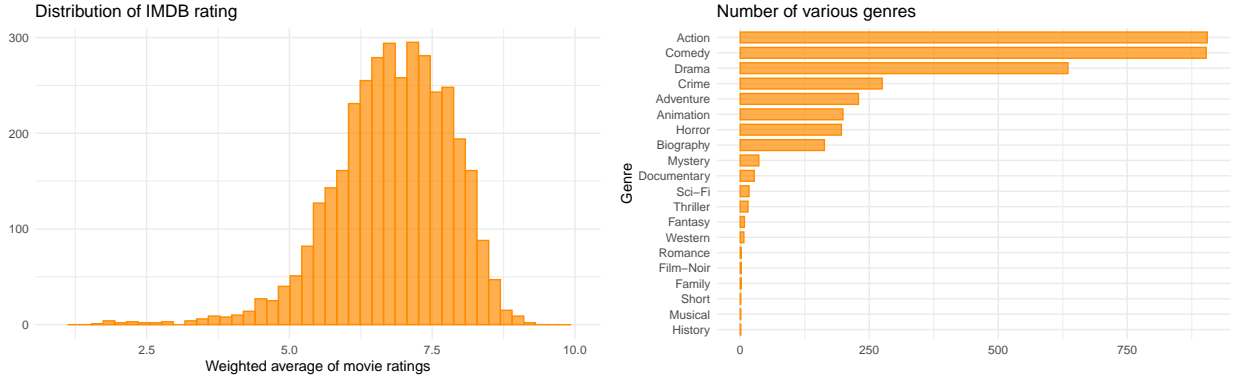


Table 1: Summary stat of IMDB ratings

Variable	mean	median	std	min	max	skew	numObs
IMDB ratings	6.8	6.9	1	1.6	9.3	-0.82	3624

The distributions (see in Appendix - Histograms) of votes, movie length, gross income and user review suggest possible log transformations. After the transformations we can see an approximately normal distribution for all except votes: it is still not close to normal distribution. In this case it is not a concern since I argued earlier why I am not intended to use this variable.

There are some **extreme values** for *user_review*, namely 8232, 6938 and 5392. But these are for films in the top 15, so it is possible, they might have a huge fan base and users left that many reviews. There is no need to drop these values. Gross income is strongly left skewed and there will be a need for log transformation in the modelling part.

The number of records drop significantly after I decided to filter based on the number of votes. The baseline was 25 thousand number of votes and my rationale behind it was the fact that when creating the top 250 list, IMDB considers only movies with minimum this vote number. I think this makes the rating more credible so I decided to apply the same. In the very end I had 3624 observations to work with.

When thinking about **representativeness** I would say that IMDB ratings are only reliable for the IMDB users who rated the movie, but not necessarily reliable indicators of what general movie audiences thought; it depends a lot on personal taste and preferences. Old movies are not that much represented since I would argue lot less registered users have seen those (even less if we consider the whole population). I also thought of the geographical coverage: only those can vote who have internet access (can register on the IMDB page) and also language might be an issue. We also cannot be 100% sure that all voters have seen the film and their vote is the true assessment of their liking. All in all this dataset captures information for registered users' preferences and the findings would indicate the favorite genre of these users only.

4. Model

4.1 Setup

As mentioned earlier, the outcome variable is rating, the weighted (on number of votes) average user rating on IMDB and the parameter of interest is genre, for which I created a new column - their number of occurrence in the dataset - for possible weighting. Before checking for potential confounder in the existing dataset, I noted several other factors which possibly can influence my outcome variable:

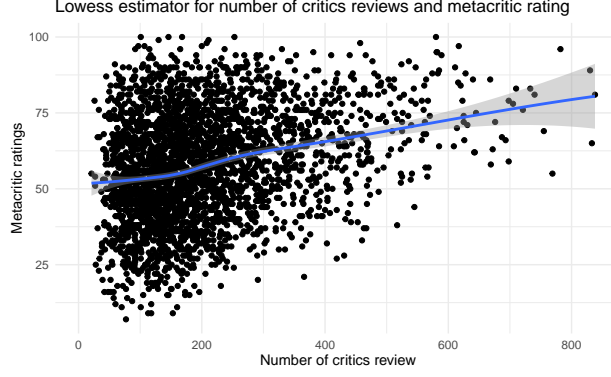
- availability of dubbing
- screening in how many countries (how worldwide the screening was) - this might relate to worldwide gross income
- lot harder to acquire and watch old films since they are not digitized or simply rare
- motion picture content rating system (whether a movie is rated only above 18 years of audience)
- leading actor/actress in the movie is famous/well known

Since I want to regress genres on ratings as the base idea, after checking these with lowess smoother (and after transforming genre into a factor), I could detect a few things:

- the most common genres are Action, Comedy, Adventure, Drama and Horror.
- Their ratings move on a wide range, while for instance Western has an average rating only above 7.5, similarly with Family or Film-Noir.



I wanted to check for a potential confounder effect between metacritic score and number of critics review since metacritic score is also based on critics reviews and scores, but there seems no pattern between the two.



In this phase I also checked correlations among my X variables. The highest correlation (where correlation is above 0.7) is between user_review and votes. This is not surprising since I filtered for observations only with at least 25 thousand votes and these movies tend to have more user reviews. ucratio is also obviously highly correlated with user_review since it is a calculated variable based on user - and critics reviews. For this reason it is enough to use either the ratio or the two other variables but not both. The moderately correlated pairs are number of votes - gross income and votes - ucratio (I would guess since it incorporates user rating which is highly correlated with ucratio).

[1] 3

Table 2: Mid - Highly correlated variable pairs

Var1	Var2	corr_val
rating	metacritic	0.7443394
metacritic	rating	0.7443394
worldwide_gross_income	votes	0.5607189
user_review	votes	0.7984952
votes	worldwide_gross_income	0.5607189
votes	user_review	0.7984952

For an example on interaction analysis see Appendix.

4.2 Modelling

The simple linear regression (1) was on ratings on the genre variable:

$$\text{Rating} = \beta_0 + \beta_1 \text{genre}.$$

At first glance on the model output, there are two categories - Horror and Fantasy - which have negative coefficients meaning these genres affect the ratings adversely. The reference category for the genre dummy is the Action category which has the highest record count in the data - **good reference category** - : the slopes doesn't seem significant, indicating no significant difference in the conditional means between the Action and other categories. For instance movies in the Adventure category received 0.4 higher rating on average than those in the Action category. I don't intend to go into more details with this model's interpretation, since its R square is really low: 0.10, similarly to models (2) and (3). So starting out with genre on the right hand side, is not very convincing. Adding duration and gross income to the model did not increase the model fit very much and the coefficients changed slightly. Hence my choice of model is the (5):

$$\text{Rating} = \beta_0 + \beta_1 \text{genre} + \beta_2 \ln(\text{duration}) + \beta_3 \ln(\text{gross_income}) + \beta_4 \text{metacritic} + \beta_5 \ln(\text{ucratic})$$

More precicely: $\text{Rating E} = 5.26 + 0.03\text{genreAdventure_} + 0.39\text{genreAnimation_} + 0.20\text{genreBiography} + 0.04\text{genreComedy} + 0.25\text{genreCrime} + 0.17\text{genreDrama} - 0.01\text{genreFamily} - 0.12\text{genreFantasy} - 0.13\text{genreHorror} + 0.09\text{genreMusical} + 0.29\text{genreMistery} + 0.52\text{genreRomance} + 0.16\text{genreSci-Fi} + 0.23\text{genreThriller} + 0.51\text{genreWestern} + 1.06\ln(\text{duration}) - 0.04\ln(\text{gross_income}) + 0.03\text{meatcritic} + 0.15\ln(\text{ucratio})$

For this model the R square jumped to 0.61, the model is accounting for 61% of the variance in the data. On the other hand the number of observations are fewer with around 800, so the size of the data is not the same as for the other previous models which performed worse. Another issue I noted is that the coefficients are quite off compared to where they were earlier: differences are around between -0.7 and +1.3. The output results suggest that on average, holding all other variables constant, a movie in category Adventure will have a higher rating wit 0.03 on average, than an Action movie, so basically almost the same. Biggest differences are with the Animation, Romance and Western categories: these genres are on average receiving higher ratings by 0.39, 0.59 and 0.51 respectively. The adverse effect of Fantasy and Horror remained but Family also turned out to have a negative coefficient. When comparing films with the same length, increasing the duration by 10 percent, we expect the rating to go up by 0.106 percent, holding all else constant. Interestingly, it seems that a 10 percent higher USD worldwide gross income we expect the rating to decrease by 0.004, but that is very low

For possible robustness check first I would focus only on movies which are categorized in one of the top 5 categories: Action, Comedy, Drama, Crime and Adventure. Another option could be of filtering for only the highly rated (let's say above 7) movies, also would try a specific decade when the films were produced (e.g. the 80s).

5. Generalization, external validity and causality

To start with, the explanatory variable *genre* is not accounting a lot for IMDB ratings alone and the model is difficult to interpret. Further control variables proved useful in explaining variation in the data, but not very convincing. It turned out that metacritic is the highest correlating variable (0.74) with rating, and when added to the model R squared increased by 0.35. If used instead of *genre* the model fit was 0.55. It is indeed an important variable, but I am not convinced that there is not a confounding variable involved - one potential variable could be *critics review* (critics review and worldwide_gross_income have a correlation of 0.4). I have listed previously some other factors I could think of influencing the movie ratings but they are hard to handle, too because they are categorical. I think the biggest con of this model is that my *genre* variable does not fully capture the information on this category. Based on this, one cannot expect to give good results if the set up is not comprehensive. The original dataset mostly contained qualitative data which also does not favor an easy model building. There was a potential right hand side variable *budget* that would make sense to incorporate: the hypothesis could be that the more money is invested in producing the film the higher it's quality and thus rating would be. Unfortunately it contained various currencies that would need to be converted to USD (I decided not to filter on observations with the USD currency because I believe that would make the analysis biased). One confounder I think could play an important role in a future analysis is whether a famous actor is playing the leading role - again, how to create the dummy can cause "headaches" (who is famous enough to be in the 'famous' category). I also believe that I have used a wide selection of data in terms of *year*, but I could not say I have a strong internal validity, since older films were also evaluated by registered users who did watch that particular film in this decade and we cannot know how a person living in the 20s would evaluate a film made around that time or a newer film. An other factor is, that ratings are highly personal and subjective and it is almost impossible to get all the angles on this matter but I would try to find solid quantitative factors indicative of these various opinions and build a model with those. From the results we can only generalize for the registered users of the IMDB webpage which is not necessarily the "single truth" for signaling how good a movie is.

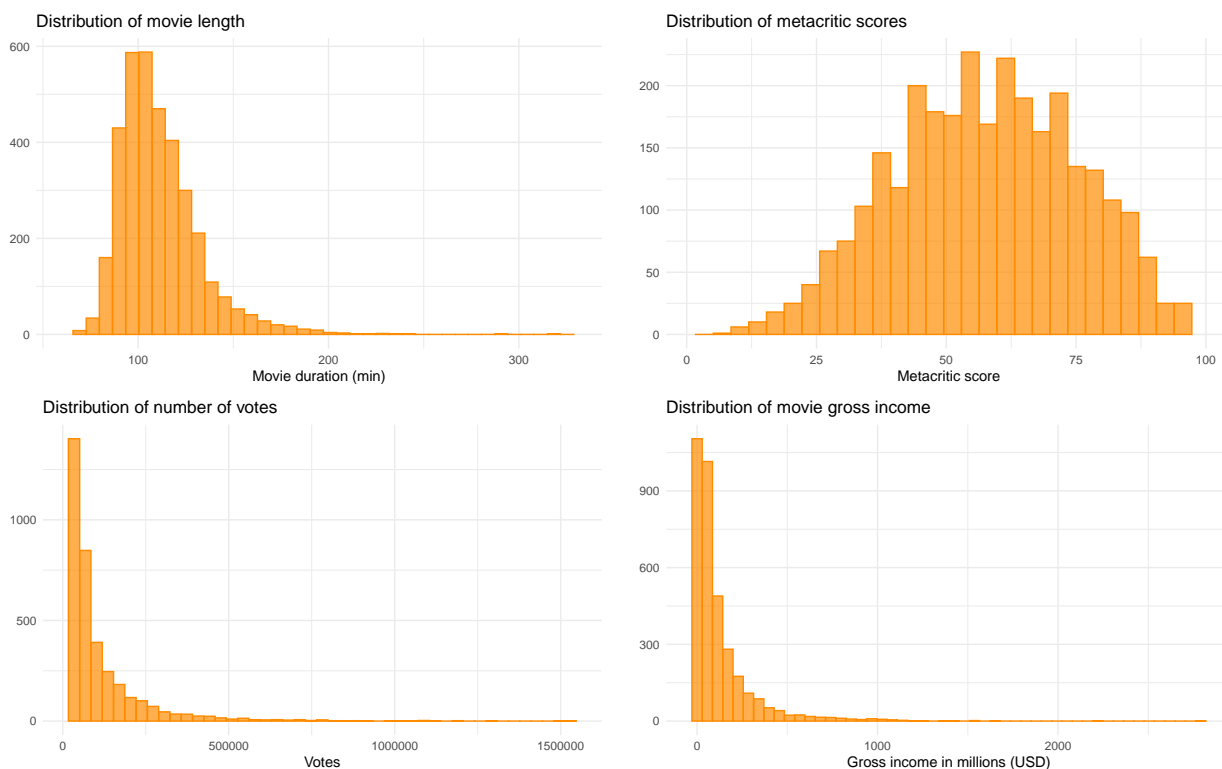
6. Summary

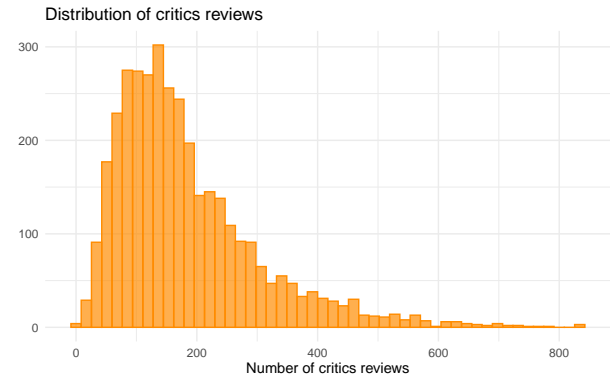
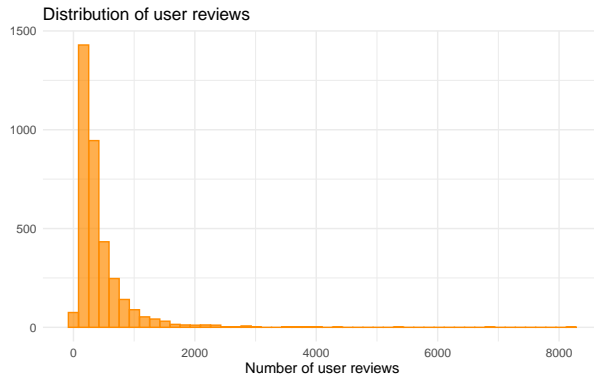
I have analyzed the relationship between average IMDB movie rating and genres. This was a simple level-level regression, which was then extended with the log transformed movie length, worldwide gross income

and user-critics ratio (number of user reviews divided by number of critics reviews) and metacritic score without transformation. With the help of the model now I know how various genres compared to the base category is associated with higher or lower ratings. I think it is informative for a producer or writer to know that topics of Family, Horror or Fantasy tend to have lower ratings (in this sample), while for instance the highest difference is for Romance: on average, holding all other variables constant, a movie in this category will have a higher rating by 0.52 on average, than an Action movie. The second highest category turned out to be Western. These might signal that despite there are very few observations with this primary genre category, these films can overtake the most common categories. This could be backed up by a robustness check analysis. The second finding was that though *metacritic* is strongly correlated with movie ratings, it does not have a strong effect on it: changing metacritic score by one, we would expect IMDB rating to change by +0.03, on average, holding all other factors constant. So could it be that users registered on IMDB are also critics who post for metacritic? Certainly.

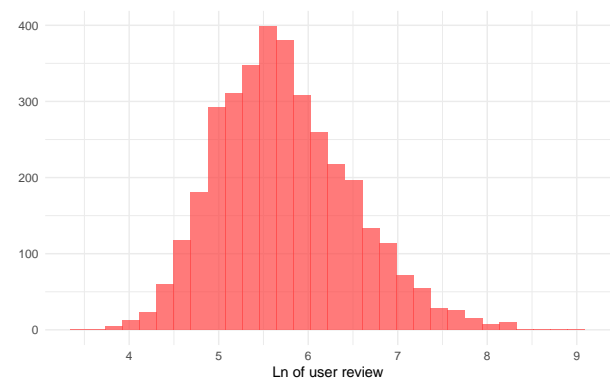
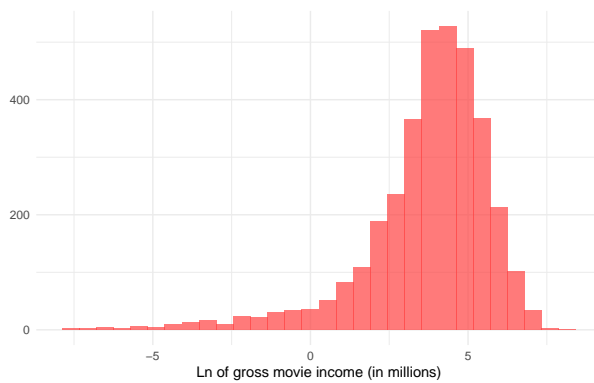
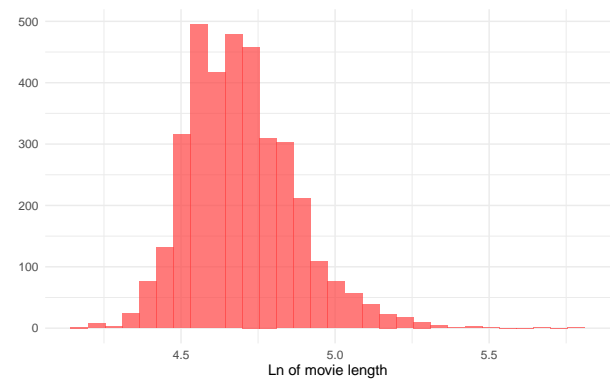
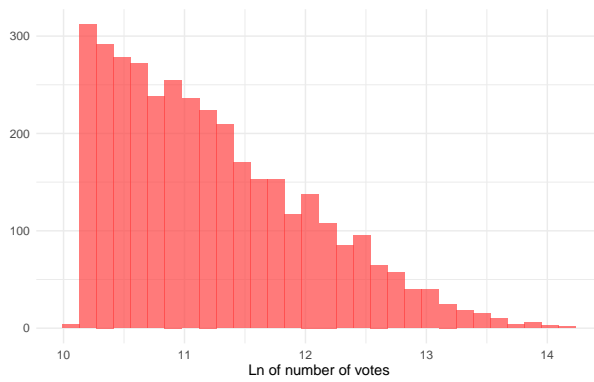
Appendices

Histograms



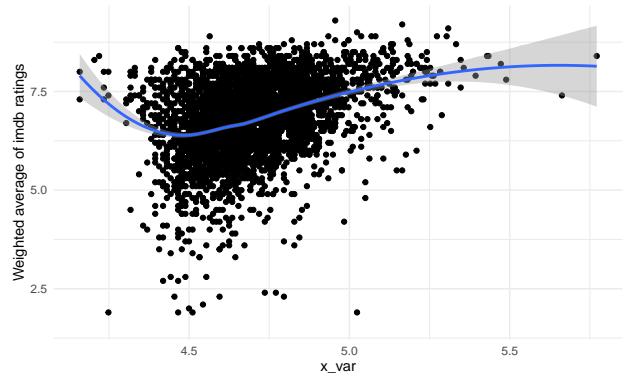
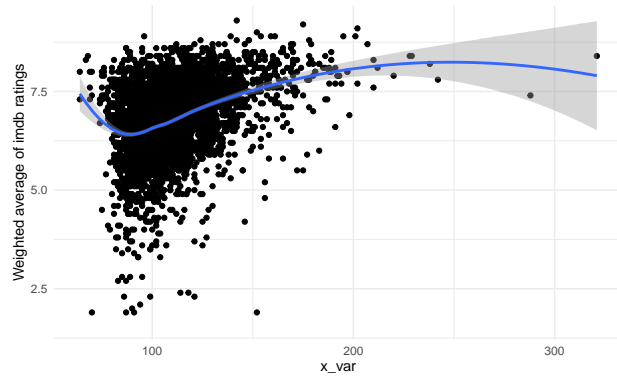


Possible log transformations

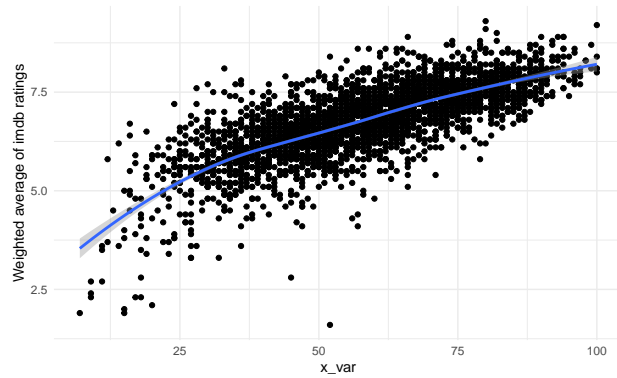


Transformations for x variables

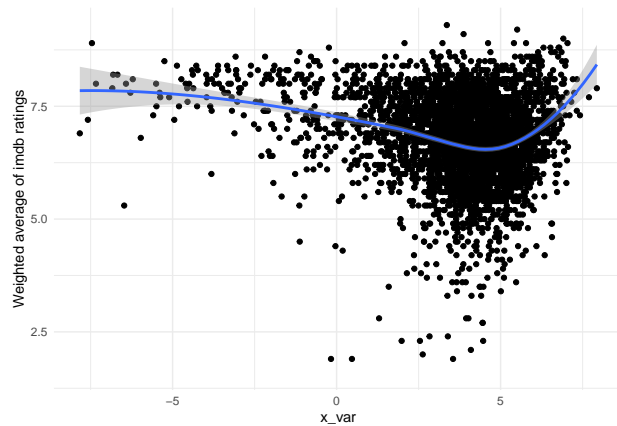
Duration Log seems better



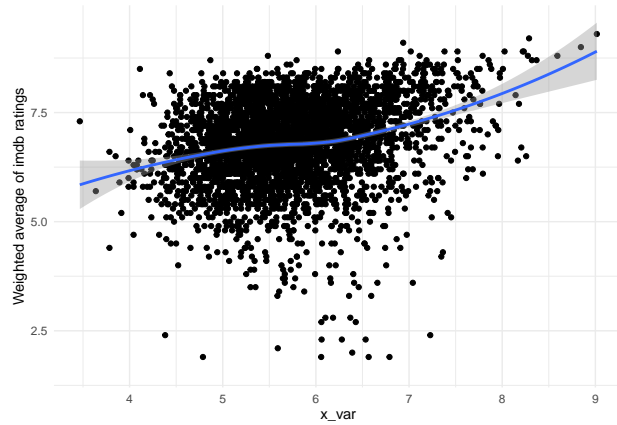
Metacritic The only variable that shows a linear pattern with IMDB rating thus this is going to be an important variable in the regression model



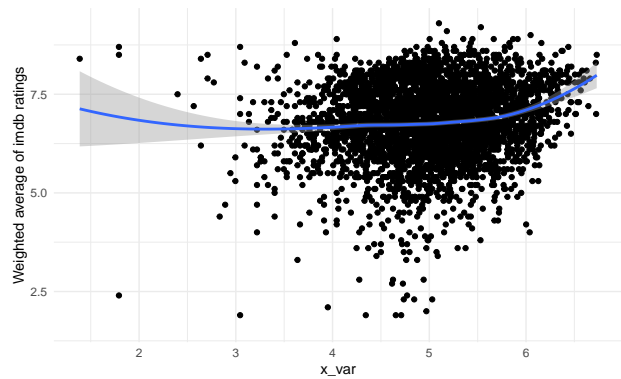
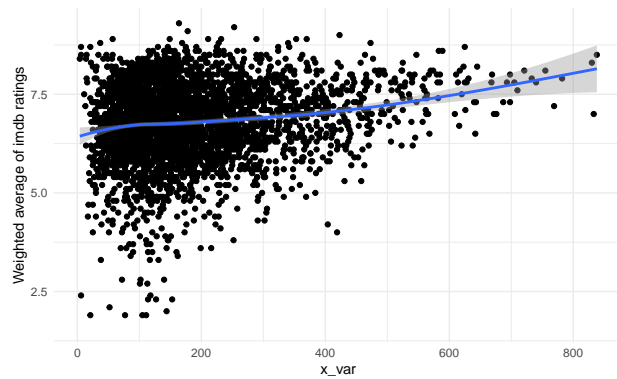
Gross income Very high heteroskedasticity, observations are centered around the middle, not very linear.



User_review Observations mainly within $\log(\text{user_review})$ 4 to 8.

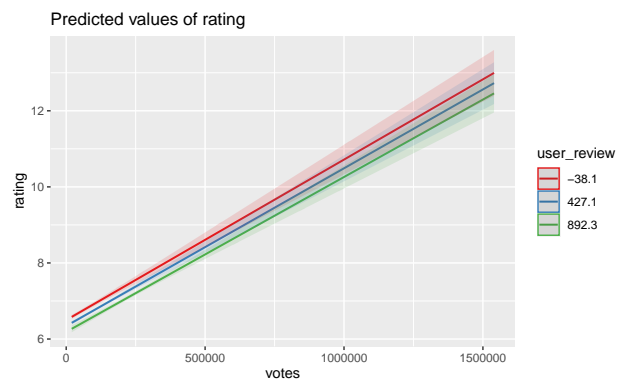


Critics review Taking the log does not really help.



Interaction

The lines between the variables votes - gross income are parallel, so no interaction occurs between them, while for ucratio - votes we can see the following graph:



The plot shows that when user review raises from low to medium there is a slight change (drop) in the output rating. The interaction term is significant at 1%, still these slope differences are not that extreme, almost parallel.