

# IMDB Rating Predictions

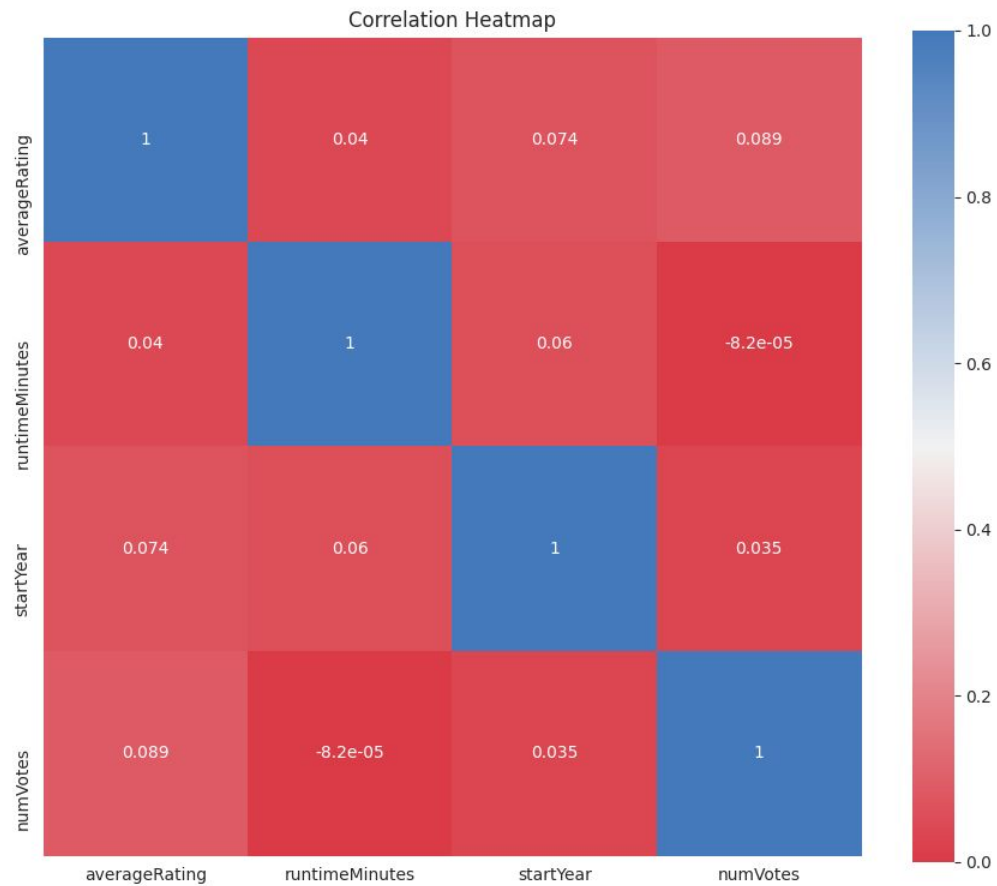
Eli Oceanak

# Goal

- Given a movie/show's basic IMDB stats other than it's average rating, I want to be able to predict the rating.
- Use things like duration, and year released to come up with a rating without having to actually watch a movie or think about it at all.

# Initial plan

- Tried to do analysis on TV shows instead of movies due to fewer people having already done analysis on them
- Datasets that contained TV shows had incomplete data, and what I ended up with had very little actual usability as predictors.

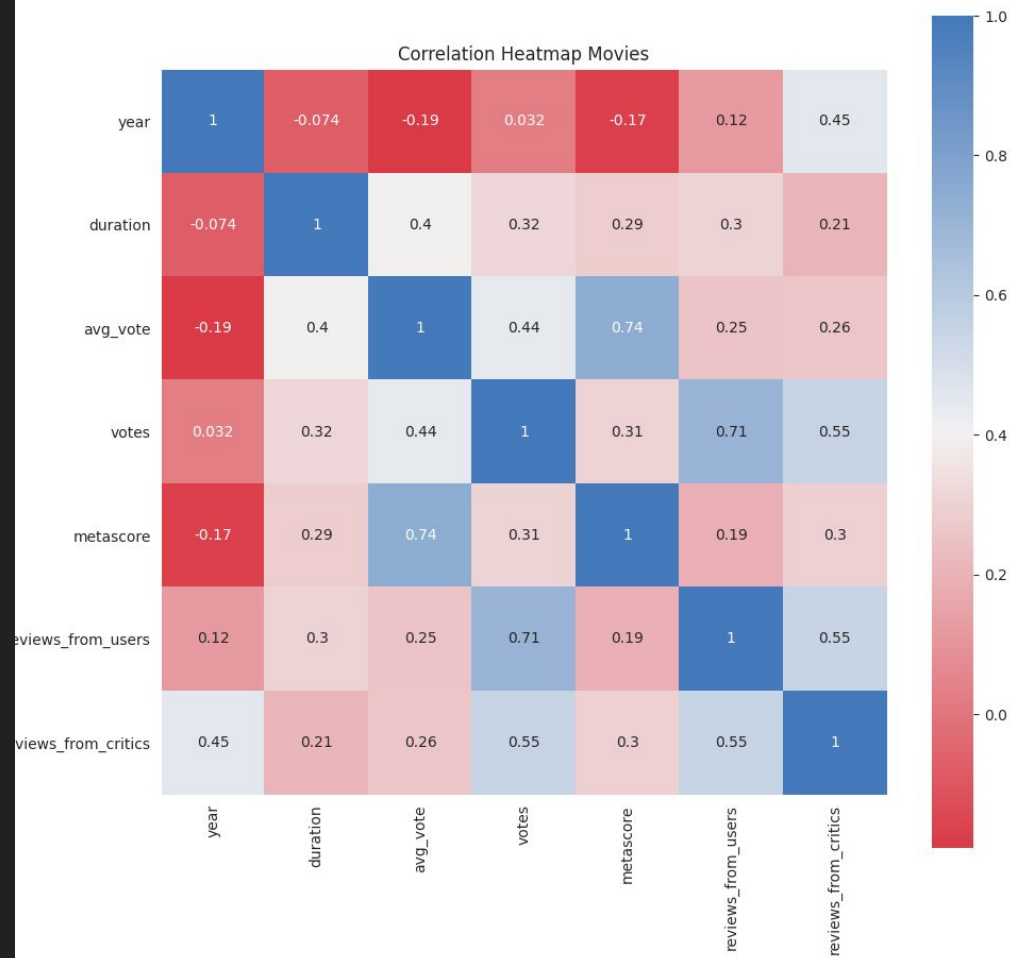


# Switched to Movies

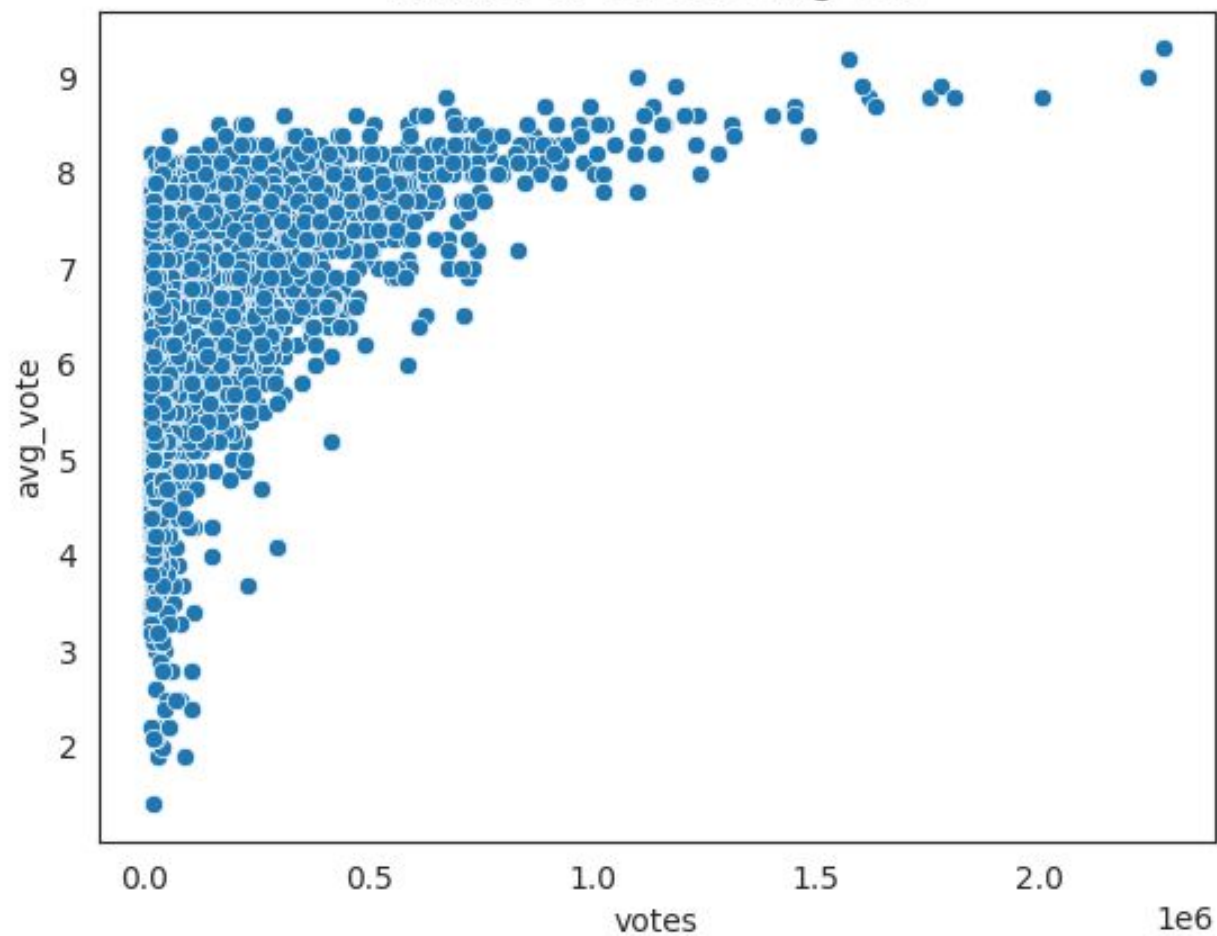
- Used this dataset -  
<https://www.kaggle.com/stefanoleone992/imdb-extensive-dataset>
- More data, more neat, and in larger quantities.
- Better predictors

# Cleaning the data

- Removed all rows with null data
- Removed all movies with  $< 10,000$  total votes cast
- Removed a few outliers with massive values that were obviously incorrect
- Ended up with just over 5,000 movies which I split 70/30 into training and testing sets.



Number of Votes vs Avg vote





# Choosing a Model

- Tried a few different models, ended up with a linear regression model.
- Well actually two different linear models
- Average error of only .65 points from just number of votes and movie duration!

```
Intercept: 4.865305933280882
Coefficients:
[('votes', 1.9001086557617272e-06), ('duration', 0.01327254030034391)]
mean squared error: 0.726093638519896
mean absolute error: 0.6474434307452884
```

```
Intercept: 3.7887129069597956
Coefficients:
[('votes', 1.0690068504436637e-06), ('duration', 0.00716504923866635), ('metascore', 0.03312036750110403)]
mean squared error: 0.3555701397621731
mean absolute error: 0.4400589079545811
```