GitHub: https://github.com/wynick27/cosi105b_weiyun

Codeclimate: https://codeclimate.com/github/wynick27/cosi105b_weiyun

Algorithm: The algorithm is to find the most similar users that watched the movie to be predicted and by the specific similarity measure and then aggregate all the results.

The similarity measure can vary, currently implemented algorithms including simple measure that returns the number of common movies, the cosine similarity measure and the pearson similarity measure. And aggregation method is either by average of all similar users' ratings or by the average of similar users' weighted by similarity, or user's average rating + average difference weighted by similarity.

The Analysis:

| Measure | Aggregation | Dataset | Mean | Stdev | RMS | Time(sec) |
|---------|-------------|---------|------|-------|------|-----------|
| common | weight_average | 20000 | 0.82 | 0.61 | 1.03 | 34.2 |
| pearson | weight_average | 20000 | 0.87 | 0.75 | 1.14 | 72.4 |
| cosine | weight_average | 20000 | 0.81 | 0.61 | 1.02 | 133.1 |
| cosine | weight_average | 2000 | 0.85 | 0.62 | 1.05 | 14.79 |
| cosine | weight_average | 200 | 0.81 | 0.60 | 1.01 | 1.70 |
| cosine | average | 200 | 0.81 | 0.60 | 1.01 | 1.69 |

From the result, we can see that pearson measurement is not good at predicting the result and cosine similarity takes too much time to compute, and the simplest similarity measure is fastest but also accurate. And the time is proportional to the dataset given the same measurement. And the aggregation method doesn't really make a difference.