

# CSE 427 final project report

## Netflix user rating prediction with collaborative filtering

Dina Elreedy, Chen Liu, Hang Yan, Hao Yan

December 14, 2016

### 1 Introduction and motivation

Recommender Systems are used in diverse applications. Online markets recommend products to customer that they may be interested in buying; social networks recommend friends or pages; online streaming services (e.g. Netflix) recommend movies to users.

Recommender systems are significant since companies can improve their sales using accurate recommender systems. Additionally, using recommender systems enhances customer satisfaction especially when he/she feels that the recommendations made by the system match their taste or needs to a great extent. Accordingly, both customers and companies would get some benefit from using recommender systems.

In this project we tackled the problem of predicting the ratings for movies from a set of users, which will guide the recommender system to recommend movies that matches his/her interest. Formally, given a set of  $n$  users and  $k$  movies and an incomplete rating matrix  $U$  of  $n$  by  $k$  entries, where  $U_{ij}$  represents the rating from user  $i$  to the movie  $j$ , our task is to predict the missing entries of  $U$ . The prediction is made possible with a large set of training data, consisting the incomplete rating of other  $m$  users for the same set of  $k$  movies.

There are multiple ways to solve this problem. The algorithm we are using is collaborative filtering, which takes the idea from  $k$ th nearest neighbor search (KNN). Given a user and the movie whose rating we want to predict, we first find similar users/movies in the training set that have the ratings on this movie, then predict the unknown rating from these known ratings. The core of the algorithm is to find similar users/movies with the testing user/movie in the training set. Given the scale of the problem, we implemented our algorithm with Spark and MapReduce and execute it on the real cloud platform.

### 2 Data analysis

The data we use is a subset of Netflix Prize data. The data consists of ratings for 17770 movies. There are ratings from 3255352 users in the training set and 100478 users in the tests set (numbers computed from the line count of TrainingRatings.txt, TestingRatings.txt and movie\_titles.txt). The size of the problem indicates that cloud computing is necessary.

For a better algorithm design, we first performed analysis to the dataset with Spark. A random subset of the data is drawn for the analysis, see Table 1 for the size of sampled data. The analysis gives us the answer of the following problems:

	Movie	User
Training set	1821	28978
Testing set	1701	27555

Table 1: The size of sampled data for analysis

## Similarity measurement

The first key design choice is which similarity measurement to use for the  $k$ th nearest neighbor search. To this end, we compute the histogram of average user ratings in the training set, see Figure 1.

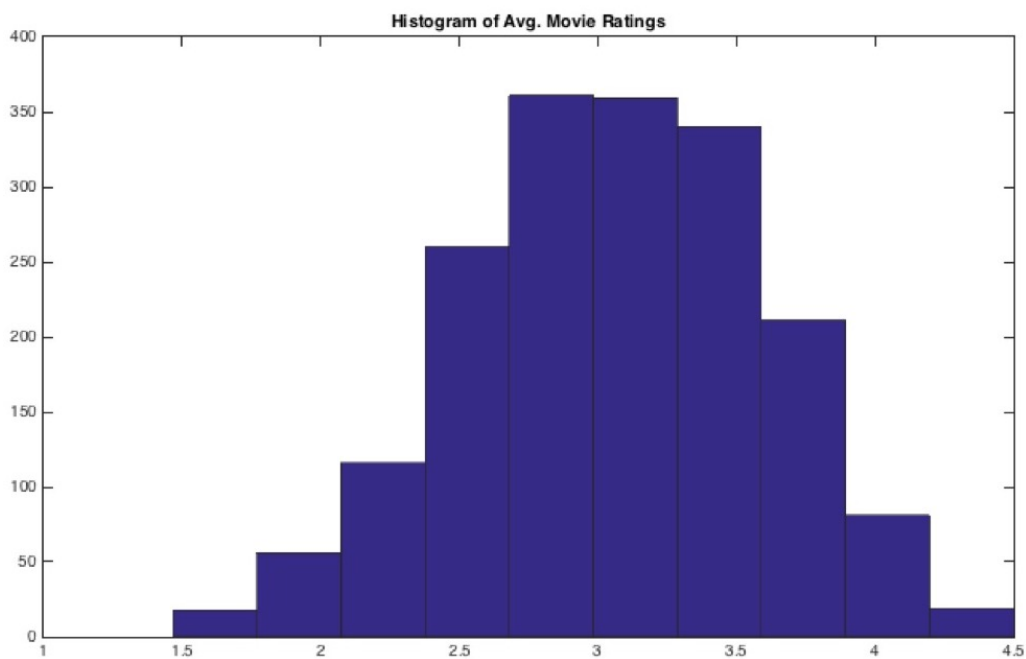


Figure 1: Histogram of average user ratings in the training set

The result suggests that most users rate normally. Therefore, we choose to use cosine similarity rather than Pearson similarity for efficiency.

## User-user similiary vs. movie-movie similarity

To fill the missing entry, we can either use user as the key and fill with ratings of other similar users, or we can use movie as the key and fill with ratings of other similar movies. The choice depends of the expected overlay of users in test and train and items in test and train. See Table 2.

The result suggests that use movie-movie similarity is more appropriate.

user-user	movie-movie
0.125467	0.439962

Table 2: Expected overly of two different keys

## Memory consumption

The algorithm will use two MapReduce jobs: the first job computes the list of users that rate each movie, and the second job computes the similarity of each pair of movies. The memory bottleneck happens at the reducer of each job. For first job, the expected memory usage is 111.5589MB and for the second job, the expected memory usage is 4.4MB.

To summarize, given a user with missing ratings, we predict those missing ratings one by one. For each movie whose rating is missing (the query movie), the algorithm first computes movies that have similar user ratings with the query movie, then predicts the missing rating with known ratings from the similar movies. The similarity is measured by Cosine similarity.

## 3 Implementation

Our algorithm consists of two major stages: the first stage is finding the most similar movies for each movie in the test set based on their overlapping ratings and the second stage is predicting the corresponding missing entries of the ratings matrix  $U$  for pairs in the test set.

For the first stage, we have implemented it using three map reduce jobs:

- Job1: This job works on the training set, it produces movie-movie pairs and their ratings by each user.

Mapper Input	$(movie_i, user_k, rate_{ik})$ triplets
Mapper Output	$(user_k, (movie_i, rate_{ik}))$
Reducer Input	$(user_k, [(movie_1, rate_{1k}), \dots, (movie_i, rate_{ik}), \dots])$ , all movies rated by user k along with their rates
Reducer Output	$(movie_i, movie_j), (rate_{jk}, rate_{jk})$ , all pairs of movies rated by user k and their rates

- Job2: This job works on the first job's output, the job mainly calculates cosine similarity between every pair of movies provided in the first job's output using their overlapped ratings.

Mapper Input	$(movie_i, movie_j), (rate_{ik}, rate_{jk})$ , rated pairs by user k.
Mapper Output	$(movie_i, movie_j), (rate_{ik}, rate_{jk})$ (It is Identity Mapper)
Reducer Input	$(movie_i, movie_j), [(rate_{i1}, rate_{j1}), (rate_{i2}, rate_{j2}), \dots, (rate_{ik}, rate_{jk})]$
Reducer Output	$(movie_i, movie_j), \cos(movie_i, movie_j)$

- Job3: After calculating similarity between each movie pairs, the third job selects the top-k similar movies for each movie. The default  $K$  we use is 10, however, we can pass  $K$  as an input parameter using command line.

Mapper Input	$((movie_i, movie_j), cos(movie_i, movie_j))$
Mapper Output	$(movie_i, (movie_j, cos(movie_i, movie_j)))$
Reducer Input	$(movie_i, [(movie_j, cos(movie_i, movie_j)) \dots])$ , list of all overlapped movies with $movie_i$ and their corresponding rates
Reducer Output	$(movie_i, [(movie_j, cos(movie_i, movie_j)) \dots])$ , Top-K similar movies per movie

For the second stage, since the top-K entries per movie is of a small size, this task does not require a Map Reduce job. Accordingly, we have implemented it using Python scripts.

## 4 Cloud Execution

The size of the problem exceeds the computation resources of a single virtual machine. We therefore use the HDInsight service on Microsoft Azure for the full execution.

We created a cluster with 2 master nodes and 4 computing nodes. Each of the computing nodes has the configuration of 4 cores CPU, 14GB RAM and 200GB local disk. See Figure 2.

The screenshot shows the Microsoft Azure portal interface for configuring a new HDInsight cluster. The 'Pricing' tab is selected, showing the 'Choose your node size' section. The cluster configuration includes:

- Cluster Name:** Enter new cluster name
- Subscription:** Microsoft Azure Sponsorship (ad8321c7)
- Cluster configuration:** Hadoop 2.7 on Linux (HDI 3.4)
- Credentials:** Configure required settings
- Data Source:** cse427storage (Central US)
- Pricing:** Please configure required settings
- Resource Group:** Create new (selected) / Use existing

The 'Choose your node size' section displays three recommended node sizes:

D3 V2 Optimized	D4 V2 Optimized	D12 V2 Optimized
4 Cores	8 Cores	4 Cores
14 GB RAM	28 GB RAM	28 GB RAM
8 Disks	16 Disks	8 Disks
200 GB Local SSD	400 GB Local SSD	200 GB Local SSD
35% faster CPU	35% faster CPU	35% faster CPU
0.62 USD/HOUR (ESTIMATED)	1.24 USD/HOUR (ESTIMATED)	0.76 USD/HOUR (ESTIMATED)

The 'Pricing' section also shows the total cost for the cluster:

Worker Nodes	Head Nodes	Total Cost
0.62 x 4 = 2.49	0.62 x 2 = 1.24	3.73 USD/HOUR (ESTIMATED)

A note indicates: 'This price estimate does not include storage costs, network egress costs, or subscription discounts. Questions? Contact billing support.'

Figure 2: Node configuration

The computing time with the cluster is listed in Table 3.

Job	Job 1 (item-item pair)	Job 2 (similarity)	Job 3 (top list)
CPU Time(s)	279.55	1516.17	13.51

Table 3: CPU time on the cluster

## 5 Results

We evaluated the performance of our algorithm and the impact of parameters on a subset of Netflix Prize dataset by the prediction accuracy and the number of unpredicted ratings in the test set. We report two accuracy measurements: the root mean square error (RMSE) and the mean absolute error (MAE).

One important parameter in our algorithm is  $K$ , the number of similar movies to consider in calculating ratings' predictions. We have tried different settings of  $K$  to investigate how it can affect the performance of our algorithm.

Figure 3 and Figure 4 shows the RMSE and MAE for our collaborative filtering algorithm using different  $K$ . The error measures, RMSE and MAE, are computed for all entries in the test set. For the entries that we could not predict (due to the small  $K$ ), we use their average ratings in the training set as our prediction. The average rating per movie is calculated by a pre-processing Map Reduce job.

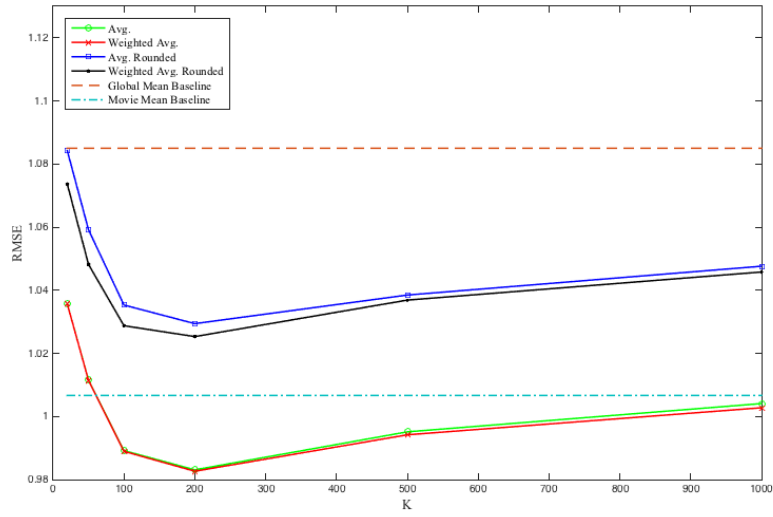
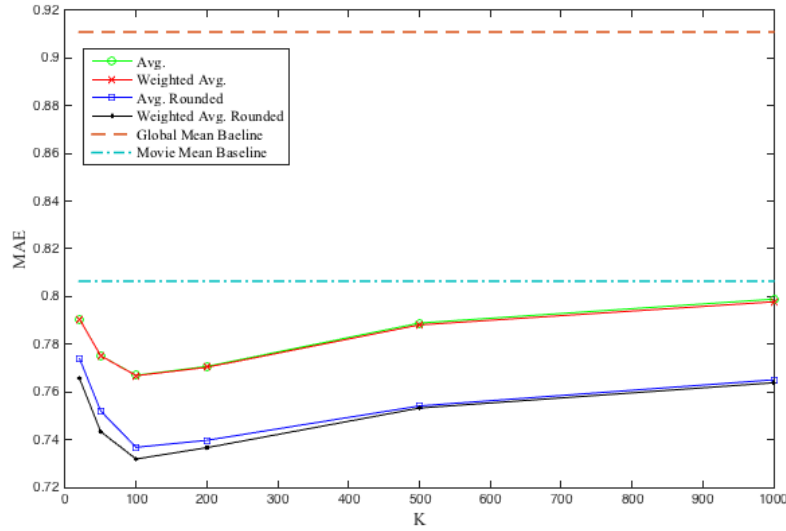
The results suggest that increasing  $K$  below a threshold helps improving the overall accuracy. However, increasing the  $K$  above that threshold will decrease the accuracy. Figure 5 shows the RMSE for predicted entries only with varying  $K$ . This suggests that the improvement of accuracy when we increase  $K$  mainly comes from the fact that we have more entries predicted. However, further increasing  $K$  over a threshold provides only insignificant gain in ratio of predicted entries, but will introduce large variance to the model, which damages the overall accuracy. We have found that the "sweet point" of  $K$  is around 200.

We also compared different methods for prediction in Figure 3 and Figure 4, including using average rating of the top  $K$  movies and using weighted (by similarity) average rating of the top  $K$  movies. We also tried rounding the predicted ratings to the nearest integers for both average rating and weighted average rating. Additionally, we compare our results to two baseline methods: the global average rating of all movies, and the average rating of each individual movie.

The results suggest that using weighted average rating is better than using unweighted average rating by only a small margin. This is due to the fact that the variance of cosine similarities between the top  $K$  movies being rated by the same user is not large. The result also shows that rounding the predicted ratings will damage the RMSE measurement, since RMSE is sensitive to noises and will amplify small errors. Therefore, we only report the unrounded results below.

Figure 6 shows the percentage of unpredicted ratings as we vary  $K$ . It shows that increasing  $K$  will reduce the percentage of unpredicted ratings, since considering more similar neighbours (movies) would increase the probability of having some of these movies rated by the testing user. However, this benefit will be insignificant when  $K$  is already large. For example, the percentage of unpredicted ratings only drops from 0.84% to 0.824% as we increase  $K$  from 500 to 1000. As discussed above, this benefit will be overwhelmed by the large variance introduced by large  $K$ .

To balance the ratio of predicted entries and the accuracy of predicted ratings, we developed an adaptive prediction scheme. We use two parameters to control the whole process: for a specific user  $i$ , and movie  $j$ , we first find  $K = 500$  candidate movies using the above algorithm. In the next step, instead of using all 500 similar movies, we only use top  $M$  movies that are rated by the user

Figure 3: RMSE for our collaborative filtering approach using different  $K$  valuesFigure 4: MAE for our collaborative filtering approach using different  $K$  values

*i*. In the case that the user  $i$  rated less than  $M$  movies among all  $K = 500$  candidates we use all candidates movies that are rated by user  $i$ .

Figure 7 and 8 show RMSE and MAE with  $K = 500$  and varying  $M$  from 5 to 100. With

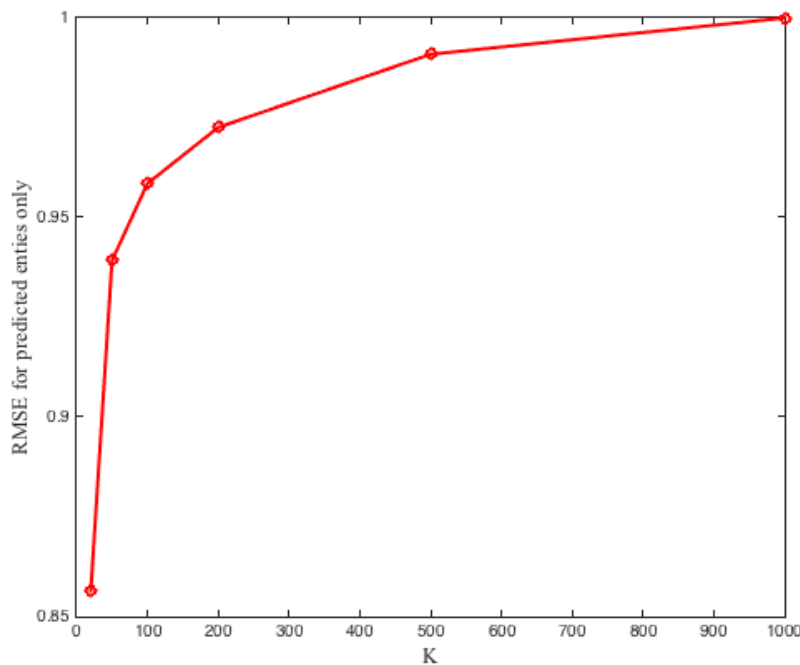


Figure 5: RMSE using weighted average prediction for predicted entries only with different  $K$  values

$M = 20$  we are able to have comparable accuracy as with  $K = 200$ , while the comparable amount of predicted entries as with  $K = 500$ , See Table 4.

Table 4 summarizes our final results for our developed collaborative filtering approach using  $K = 500$  and  $M = 20$ . Results are presented in terms of prediction accuracy using RMSE and MAE, and for the percentage of unpredicted ratings. The table also shows the comparison against two baselines methods: the global mean rating and movie mean rating. Our collaborative filtering approach have better accuracy than the two baselines with low percentage of unpredicted entries.

	RMSE	MAE	Percentage of unpredicted ratings
Collaborative Filtering	<b>0.98099</b>	<b>0.76806</b>	0.84%
Global Mean Baseline	1.08502	0.91094	-
Movie Mean Baseline	1.0067	0.8065	-

Table 4: Result of our collaborative filtering and the comparison against two baseline methods.

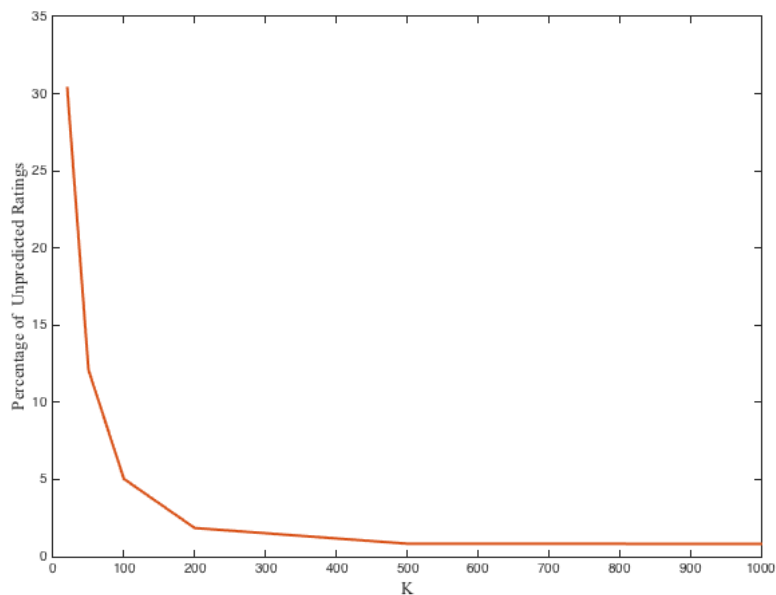


Figure 6: Percentage of unpredicted ratings using different  $K$  values

## 6 Conclusion

In this project, we have developed a system that predicts missing ratings of a movie for a user using collaborative filtering approach which helps in recommendation systems as these ratings can guide the company like Netflix which movies to recommend to which users. First, we have analyzed the data using Spark and designed the algorithm based on this analysis results. Then, the collaborative filtering algorithm is implemented with MapReduce framework and executed on Microsoft Azure Hadoop service. We have also improved the original collaborative filtering and the prediction algorithm to balance the prediction accuracy and number of predicted entries. Experimental results demonstrate the effectiveness of our system and show the impact of several parameters.



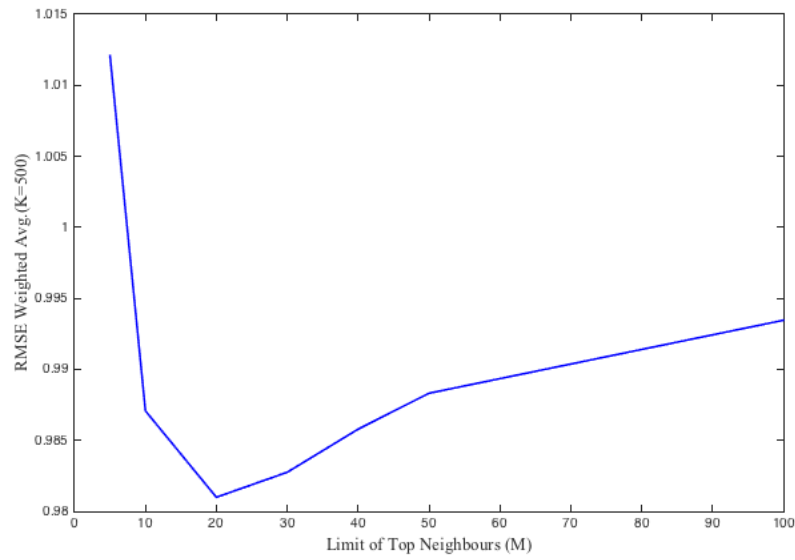


Figure 7: RMSE with varying  $M$  using weighted average prediction and  $K = 500$ .

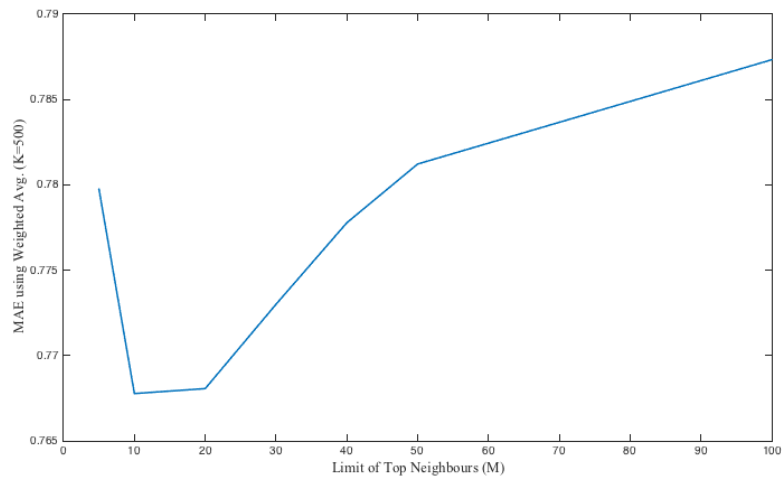


Figure 8: MAE with varying  $M$  using weighted average prediction and  $K = 500$ .

## Appendix A: group collaboration

Our group consists of four members. The tasks for each group member is listed in Table 5:

Name	Task
Dine Elreedy	Algorithm design, collaboarative filtering, report
Chen Liu	Algorithm design, collaboarative filtering, prediction
Hang Yan	Algorithm design, cloud platform, report
Hao Yan	Algorithm design, data analysis

Table 5: Tasks for group members