# CS 6001 Applied Spatial and Temporal Data Analysis - Spring 2017 - Homework 4

Dalton Cole

April 3, 2017

## 1  Introduction

The purpose of this assignment is to become familiar with recommender systems. In the video that Dr.Fu asked us to watch, User-based collaborative filtering and item-based collaborative filtering was discussed. This assignment gives us first hand experience at these two filtering methods. User-based collaborative filtering is done by finding the similarity between different users based on the items, in the video's case, Netflix shows watched. Item-based collaborative filtering is essentially the transpose of user-based. Similarity between items are used instead of between users. In most cases, this provides more reliable data. This is in part because, normally, there are fewer catalog items than there are users. This reduces the dimensionality of the data.

These two methods are variations of the kNN algorithm. The basic formula is:

$$\hat{r}_{ui} = \frac{\sum\limits_{v \in N_i^k(u)} \text{sim}(u, v) \cdot r_{vi}}{\sum\limits_{v \in N_i^k(u)} \text{sim}(u, v)}$$

Three more algorithms are used: SVD, PMF, and NMF. The SVD algorithm was popularized by Simon Funk, as discussed in the KNN tutorial video. It uses a prediction value:

$$\hat{r}_{ui} = \mu + b_u + b_i + q_i^T p_u$$

The following regularized squared error should be minimized using stochastic gradient descent:

$$\sum_{r_{ui} \in R_{train}} (r_{ui} - \hat{r}_{ui})^2 + \lambda \left(b_i^2 + b_u^2 + ||q_i||^2 + ||p_u||^2\right)$$

PMF is the unbiased version of SVD:

$$\hat{r}_{ui} = q_i^T p_u$$

NMF is very similar to SVD as well. The difference between the two is that the user and item factors are kept positive in NMF.

MAE and RMSE are two matrices that are used to evaluate performance. Mean absolute error (MAE) can be seen in Figure 1. Root-mean-square deviation (RMSD) can be found in Figure 2.

$$\mathrm{MAE} = \frac{1}{n}\sum_{i=1}^{n}|f_i - y_i| = \frac{1}{n}\sum_{i=1}^{n}|e_i|$$

Figure 1: Mean Absolute Error Formula

$$\mathrm{RMSD} = \sqrt{\frac{\sum_{t=1}^{n}(\hat{y}_t - y_t)^2}{n}}$$

Figure 2: Root-Mean-Square Deviation Formula

## 2  Experiment

The code used for this segment can be found in *HW4.py*. The data used in this experiment is located in *restaurant_ratings.txt*, which is a text file given by Dr.Fu for this assignment. The goal of this experiment is to see which algorithm performs best given this data set. Three-fold cross validation was used. To keep each fold consistent across each algorithm and across seperate tests, the random seed generator was seeded with the value *0*.

## 3  Results

The results of this experiment can be seen in Figures 3, 4, 5, and 6. As can be seen in Figure 3, SVD out performed the other methods. In this case, item-based collaborative filtering barely out performed user-based collaborative filtering. Figure 4 shows how important a similarity metric can be. Using the MSD similarity metric, item-based out performs user-based collaborative filtering. Using cosine or Pearson however gives user-based the edge.

Figure 5 shows that when k = 25, 27, and 28, RMSE was as its minimal value for user-based collaborative filtering. This implies that looking at the mid 20s to higher 20s neighbors yields the best results for. When k = 43 gave item-based collaborative filtering its best RMSE value. Item-based out performed user-based when $k \geq 30$. When $k < 30$, user-based out performs item-based.

## 4  Conclusion

Different machine learning algorithms produce different results. Different error measuring formulas also produce different results. It is important to keep in mind both of these facts when selecting algorithms. It is important to be consistent with the error measurements. In this report, RMSE was primarily used to compare algorithms. It has been shown, that with the given data set, and with the given random folds, SVD is the
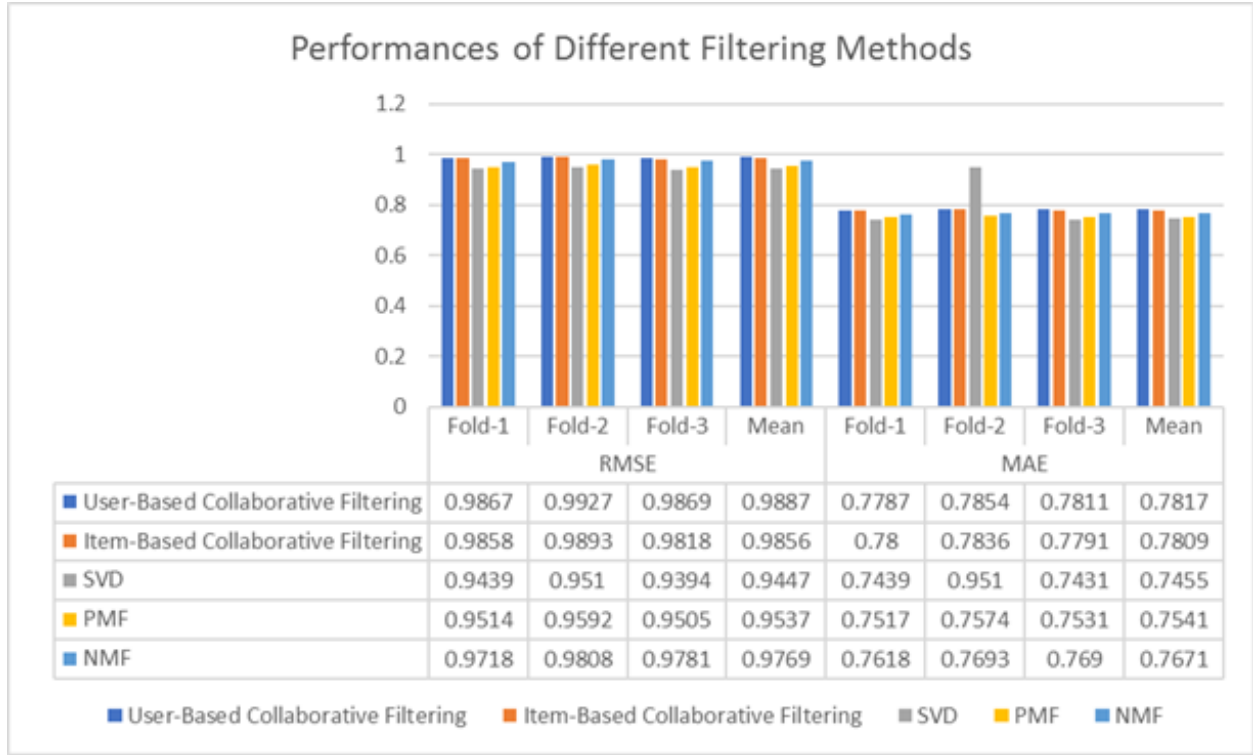
Figure 3: Performance of Different Filtering Methods

| | Fold-1 | Fold-2 | Fold-3 | Mean | Fold-1 | Fold-2 | Fold-3 | Mean |
|---|---|---|---|---|---|---|---|---|
| | | RMSE | | | | MAE | | |
| User-Based Collaborative Filtering | 0.9867 | 0.9927 | 0.9869 | 0.9887 | 0.7787 | 0.7854 | 0.7811 | 0.7817 |
| Item-Based Collaborative Filtering | 0.9858 | 0.9893 | 0.9818 | 0.9856 | 0.78 | 0.7836 | 0.7791 | 0.7809 |
| SVD | 0.9439 | 0.951 | 0.9394 | 0.9447 | 0.7439 | 0.951 | 0.7431 | 0.7455 |
| PMF | 0.9514 | 0.9592 | 0.9505 | 0.9537 | 0.7517 | 0.7574 | 0.7531 | 0.7541 |
| NMF | 0.9718 | 0.9808 | 0.9781 | 0.9769 | 0.7618 | 0.7693 | 0.769 | 0.7671 |



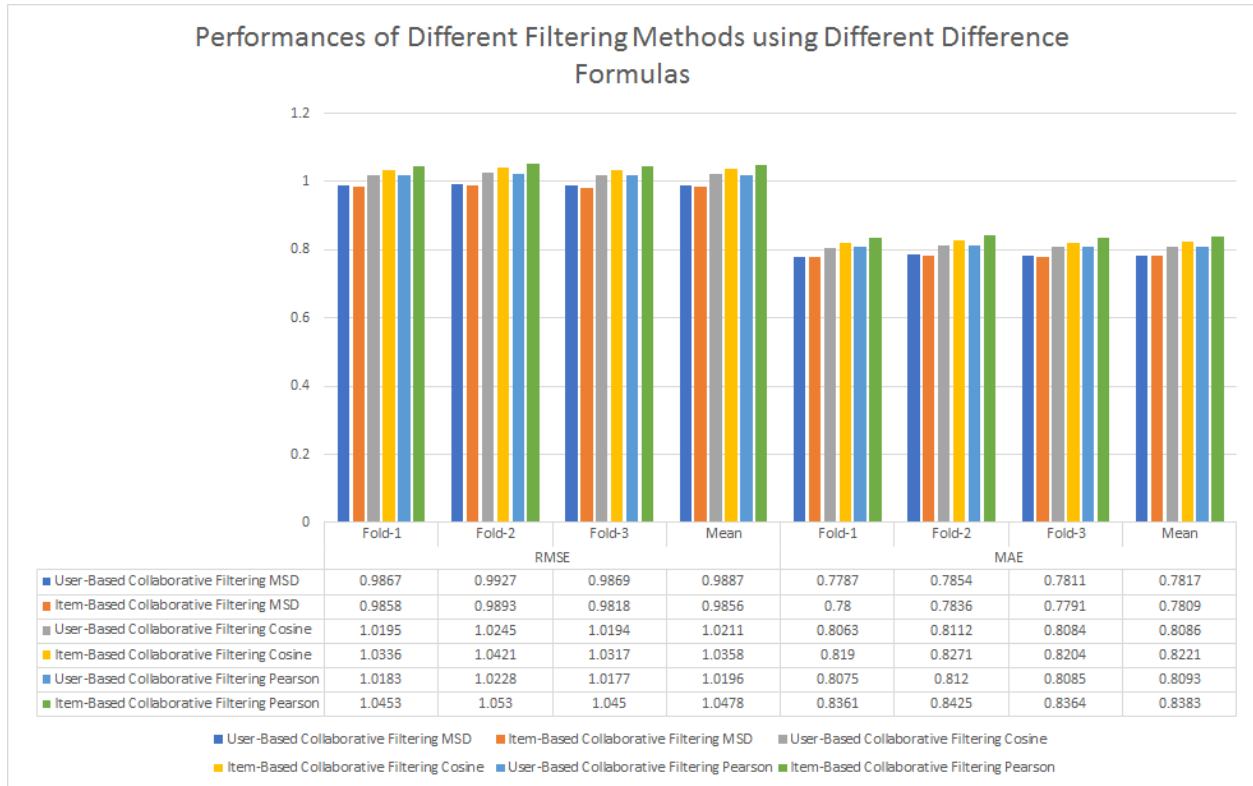| | Fold-1 | Fold-2 | Fold-3 | Mean | Fold-1 | Fold-2 | Fold-3 | Mean |
|---|---|---|---|---|---|---|---|---|
| | | RMSE | | | | MAE | | |
| User-Based Collaborative Filtering MSD | 0.9867 | 0.9927 | 0.9869 | 0.9887 | 0.7787 | 0.7854 | 0.7811 | 0.7817 |
| Item-Based Collaborative Filtering MSD | 0.9858 | 0.9893 | 0.9818 | 0.9856 | 0.78 | 0.7836 | 0.7791 | 0.7809 |
| User-Based Collaborative Filtering Cosine | 1.0195 | 1.0245 | 1.0194 | 1.0211 | 0.8063 | 0.8112 | 0.8084 | 0.8086 |
| Item-Based Collaborative Filtering Cosine | 1.0336 | 1.0421 | 1.0317 | 1.0358 | 0.819 | 0.8271 | 0.8204 | 0.8221 |
| User-Based Collaborative Filtering Pearson | 1.0183 | 1.0228 | 1.0177 | 1.0196 | 0.8075 | 0.812 | 0.8085 | 0.8093 |
| Item-Based Collaborative Filtering Pearson | 1.0453 | 1.053 | 1.045 | 1.0478 | 0.8361 | 0.8425 | 0.8364 | 0.8383 |

Figure 4: Performance of Different Filtering Methods using Different Difference Formulas
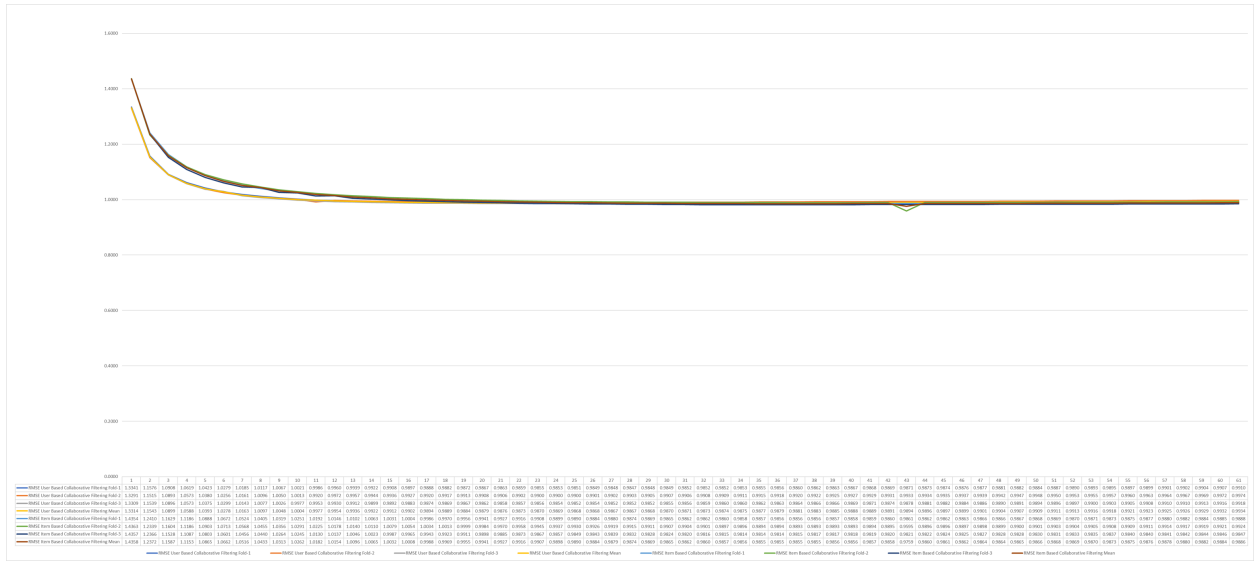
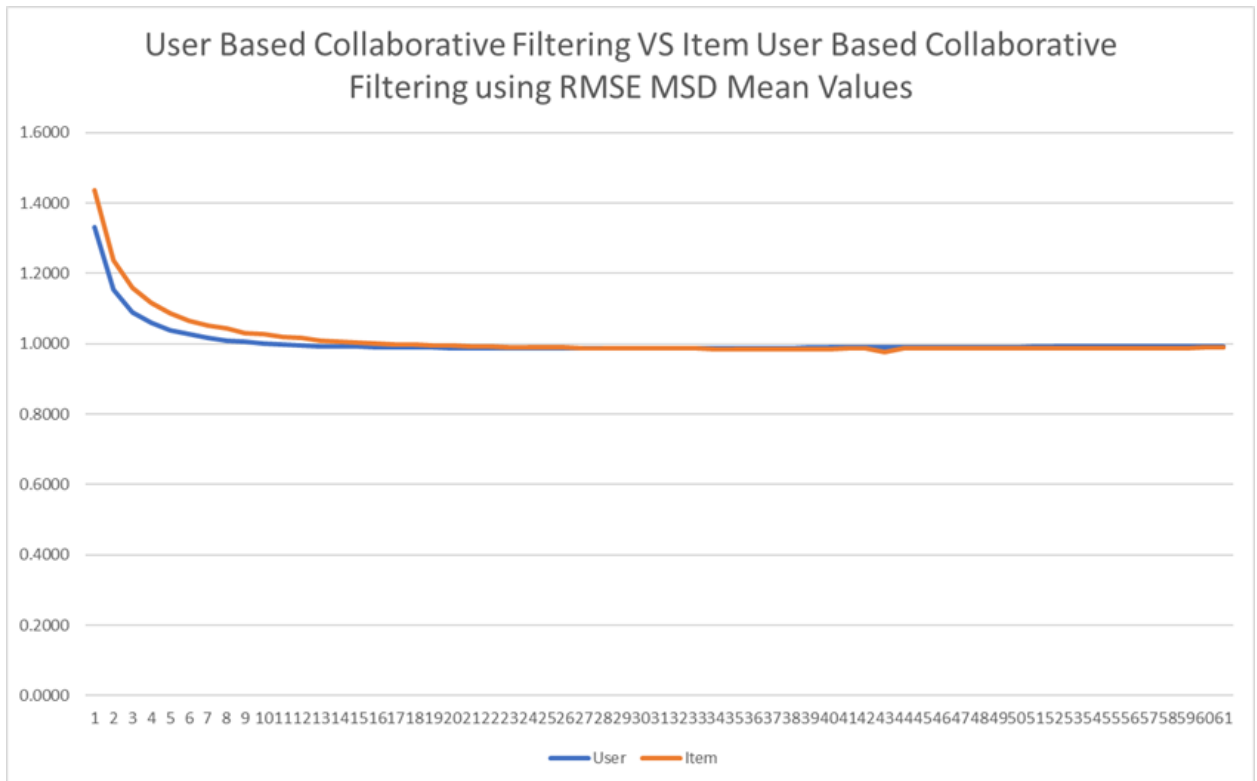Figure 5: RMSE Filtering with each Fold and Mean of the Three Folds



Figure 6: User-Based Collaborative Filtering VS Item-Based Collaborative Filtering using RMSE MSD Mean Values

best filtering algorithm out of the five in the experiment. Between user-based and item-based collaborative filtering, item-based generally out performed user-based. Since these two methods are variations of the kNN method, it is important to keep in mind the value of k. When k has a value less than 30, user-based collaborative filtering out performed item-based. This will vary depending on the data set, but the general trend that user-base will out perform item-based up until a certain threshold is expected.

Different similarity measurements matter highly as well. As shown in Figure 4, item-based Pearson method performed the worst, however, out of the user-based methods, using the cosine similarity performed the worst. Different similarity measurements vary on different algorithms and data sets.