

Movie Recommendation

Arighna Roy

I. Definition

Project Overview

A recommendation engine can be used for a lot of applications such as categorization of products and responding to users' individual interest. In this project, I have tried to recommend products (movies) to users (viewers) based on the viewer to movie rating data. The underlying technique to achieve the goal is collaborative filtering. I have used both user based and content based collaborative filtering and made a comparative study.

There are multiple techniques to build a recommendation engine such as item set mining, collaborative filtering, content based filtering and knowledge based filtering. Item set mining is used in filtering out the items which are frequently used together. Collaborative filtering can be used when users' ratings for different products are available. For Content based and knowledge based filtering, we need information related to products or some prior knowledge about the market.

Problem Statement

I have a set of users for which the ratings to different movies are available. Now, for a new user with his ratings for few of the existing (rated by existing users) movies is taken as an input. I aim to predict the rating of that user to a movie which he has not rated, but that movie is rated by other users. Let me explain the problem with an example. Suppose, there are 3 users a,b,c. a rated movie 1,2,3; b rated movie 2,3,4; c rated 3,4,5. Now a new user d rated movie 1,4,5; I aim to predict the rating of movie 2,3 for user d.

Step1: preprocess the user rating file. We removed the outliers from the rating dataset. Users who rated too low or too high number of movies are considered as outliers. Similarly, movies rated by too low or too high number of users are considered as outliers.

step2: generate a python dictionary for the user to movie ratings. Each key of the dictionary is a user id(suppose u). Values of each key(u) is another dictionary. The key of each dictionary is a movie id (suppose m) and the corresponding value is the rating of that movie(m) by the user (u).

step3: build a similarity matrix between users. Each key of the dictionary is a user id(suppose u1). Values of each key(u1) is another dictionary. The key of each dictionary is another user id (suppose u2) and the corresponding value is the similarity between the two users (u1 and u2).

step4: build a similarity matrix between movies. Each key of the dictionary is a movie id(suppose m1). Values of each key(m1) is another dictionary. The key of each dictionary is another movie id (suppose m2) and the corresponding value is the similarity between the two movies (m1 and m2).

step5: loop through each user-movie pair and find the predicted rating of the movie for that user using both user based and movie based collaborative filtering.

step6: Compare the predicted value and actual rating. Calculate the MSE.

Metrics

I use the same user-movie pairs from the training set and compare it with the existing rating. Suppose, user A has rated x (<10) to movie 1. I predict the ratings of movie 1 for the user a using my method. The rating of movie 1 for user A is ignored. And the performance of the prediction is measured based on the comparison with the value x (actual rating).

There are many performance metrics which can be considered to evaluate the correctness of the algorithm.

Root Mean Squared Error (RMSE), Mean Absolute Error (MAE), Mean Signed Deviation (MSD) are the most widely used error metrics.

I chose **RMSE** over **MAE** and **MSD**. **MAE** is the absolute deviation of the predicted value. **RAE** is good where I want the impact of deviation to be equal at every scale as opposed to in **RMSE** the effect of deviation exponentially increases. In my opinion, **RMSE** is more appropriate here because the scale is small here and the number of rating option is also small (5). **MSD** is not applicable here because we don't need to know whether the predicted value is higher or lower than the actual value. We just want to find the absolute value.

To evaluate the performance of the algorithm, I have used **root-mean-square error** as the performance metric.

$$RMSE_{Errors} = \sqrt{\frac{\sum_{i=1}^n (\hat{y}_i - y_i)^2}{n}}$$

To construct the r.m.s. error, you first need to determine the residuals. Residuals are the difference between the actual values and the predicted values. I denoted them by $\hat{y}_i - y_i$,

where y_i is the observed value for the i th observation and \hat{y}_i is the predicted value.

II. Analysis

Data Exploration

The data set contains a rating file. I have splitted the data into training and testing dataset. Below is a sample of 10 records from the rating dataset:

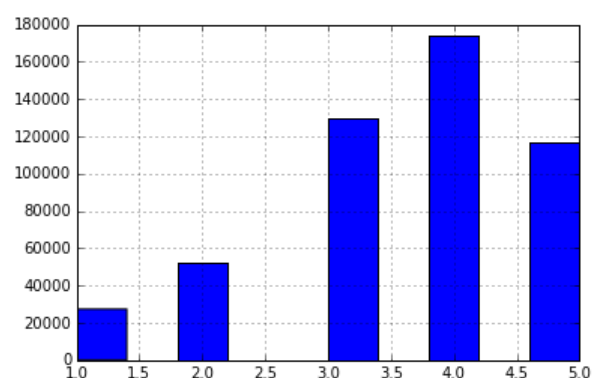
	user	movie	rating	id
0	2783	1253	5	2783_1253
1	2783	589	5	2783_589
2	2783	1270	4	2783_1270
3	2783	1274	4	2783_1274
4	2783	741	5	2783_741
5	2783	750	5	2783_750
6	2783	924	5	2783_924
7	2783	2407	4	2783_2407
8	2783	3070	3	2783_3070
9	2783	208	1	2783_208

Below is a short summary of the dataset:

Total number of ratings:	500100
Total number of users:	3255
Total number of movies:	3551

Below is a summary statistics and the histogram of the ratings:

count	500100
mean	3.60222
std	1.11469
min	1
25%	3
50%	4
75%	4
max	5



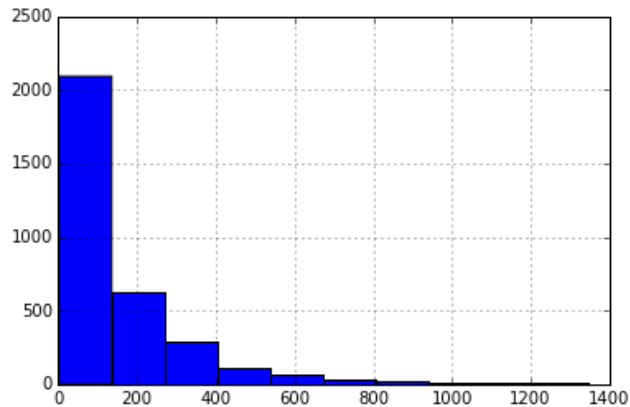
Exploratory visualization

Below are 5 users who rated the highest number of movies and the count of the movies they rated:

User_id	number of rated movies
3618	1344
5795	1272
4344	1271
4510	1240
4227	1222

Below is short summary of the count of the rated movies by each user and its histogram:

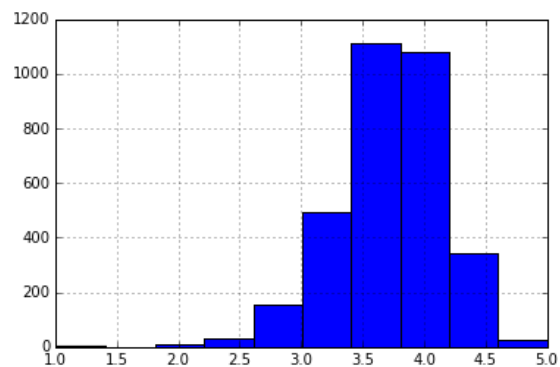
count	3255
mean	153.640553
std	173.90484
min	2
0.25	42
0.5	92
0.75	197
max	1344



We can clearly see that most of the users rated very few number of movies, where as very few rated more than 400 movies.

Below is short summary of the average ratings of each user and its histogram:

count	3255
mean	3.713232
std	0.431985
min	1.015385
25%	3.455723
50%	3.752475
75%	4
max	5

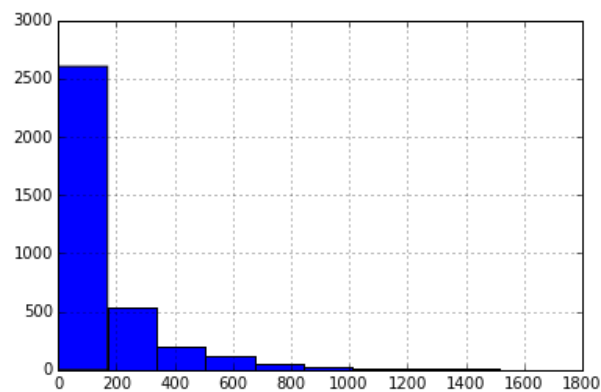


Below are 5 movies which are rated the highest number of time and the count of the usersrated those:

Movie_id	number of users rated
2858	1684
1196	1585
260	1573
1210	1539
127000%	1396

Below is short summary of the count of the users rated each movie and its histogram:

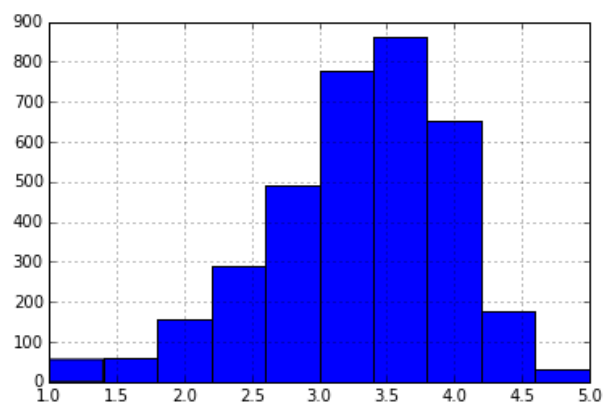
count	3551
mean	140.833568
std	200.241913
min	1
25%	18
50%	66
75%	181
max	1684



We can clearly see that most of the movies are rated by very few number of users, where as very few movies are rated by more than 400 users.

Below is short summary of the average ratings of each movie and its histogram:

count	3551
mean	3.283631
std	0.702784
min	1
25%	2.869126
50%	3.375427
75%	3.792778
max	5



Algorithms and Techniques

I have used the collaborative filtering (both user based and movie based) to find the nearest neighbors for predicting the rating. For each pair of user and movie in the testing dataset, I have omitted the rating and made a prediction based on the all other ratings. Then compared the prediction.

Benchmark

The winner of the Netflix competition for movie rating prediction achieved RMSE=0.86. They followed such a complicated method which took almost 3 years to complete. For a simple and standard approach like this, I would keep a target of an $RMSE \leq 1$ which is practical. One point to note that the difference in RMSE might look very small but in a scale of 5, the difference in RMSE of 0.14 is actually pretty large.

III. Methodology

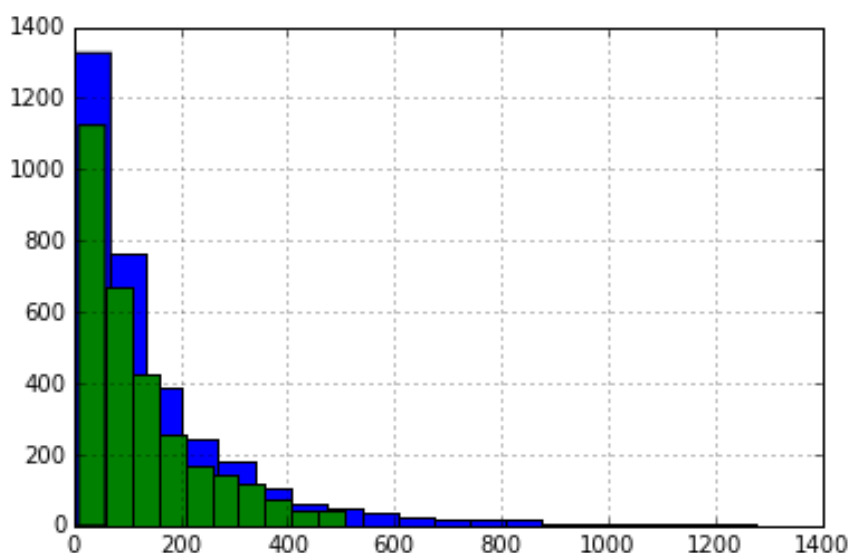
Data Preprocessing

It's always a good idea to pre-process the data before starting the actual analysis or building the application. This dataset doesn't have any missing values or non-numeric values. I have done outlier detection and removal as part of pre-processing.

1% of the users rated less than 9 movies. 9 is a very small number in terms of movies rated. It's very difficult to judge someone's taste of movie based on the ratings of just 9 movies. So remove the users who rated less than 9 movies for my analysis. 5% of the users rated more than 507 movies. I consider them someone having variety of taste in the choice of movies. Those users should not be used for the analysis. The reason we choose 5% in the upper tail and 1% in the lower tail is because we have a right skewed distribution for the count of movies rated by each user.

I want to consider the mean and the standard deviation of the distribution while finding the outliers in the upper tail of the distribution. 77 users rated number of movies which are 3 standard deviation above the mean of the distribution. But I am dealing with the highly skewed distribution. I would rather prefer the quantile data than the variance. Also, I set these values after a careful consideration of the count of movies.

Below is the histogram of the count of movies rated by the each user before (blue) and after (green) the removal of the outliers.



Implementation

The implementation of my project is broken down into few ipython notebooks.

Visualization:

In this notebook, I have performed the statistical analysis of the data and further exploratory analysis.

I have used pandas and numpy to read the data and visualize the summary statistics of the data.

Pre-processing:

In this notebook I have detected the outliers and prepared the final data for analysis.

Reco:

In this notebook, I have implemented the prediction algorithm

I have implemented user based prediction, movie based prediction and the average of the previous 2 values to test the accuracy of the result.

Results:

Here, I have generated the plots based the prediction of movie rating.

I have used Collaborative filtering for predicting the rating of a movie by one user. This is also called the social filtering. Let me explain this idea with an example. Suppose, I hang out with some friends on regular basis. The group is quite diverse based on their tastes of food. Few of them have similar tastes as mine. I figured that out when we went to restaurants together and they liked or disliked the food along with me. Now, I want to go to a restaurant alone and want review of the restaurant. More than the online reviews, I would like to have the reviews from those friends who have similar tastes of food as mine. I can say that the rating I am going to give to a restaurant is likely to be close to the ratings those people have already given. This is exactly the same idea behind user based collaborative filtering. In our problem, the restaurant is replaced by movie and my friends are the online users. I have used user based and movie based collaborative filtering for finding out the rating of a movie by one user.

Outline of the implementation of collaborative filtering

Let me discuss the outline of implementation of prediction algorithm here.

Predict rating by **user** based collaborative filtering:

- Build the similarity matrix for users
 - For a pair of users, find the movies which both the users have rated. Now find the distance between the users' ratings for the common movies. Use one of the distance metrics to calculate this.
 - Repeat this for each user pair of users
- Pick a combination of user and movie from the testing set
- Filter the users which have rated that movie
- Find the users with min distances from the chosen user
- Now find the average of the ratings given by those users to that movie

Predict rating by **movie** based collaborative filtering:

- Build the similarity matrix for movies with a similar process followed for user similarity matrix
- Pick a combination of user and movie from the testing set
- Filter the movies which rated by that user
- Find the movies with min distances from the chosen movie
- Now find the average of the ratings given by that user to those movies

Refinement

While approaching this problem, I thought 2 methods to consider: collaborative filtering and content based filtering. The major dataset I need for collaborative filtering is the ratings history and for content based filtering is the movie features. But the dataset I considered had only few features of each movie which were not sufficient to make a good prediction about the rating. So I chose collaborative filtering.

Now, there are two different methods to consider in collaborative filtering, user based rating and movie based rating. I decided to perform both and do a comparison study.

I first included all the data in training to predict the rating of each movie by each user. This wasn't such a bad idea because K-nearest neighbour is a lazy learner. But to make the result beyond any confusion, I split the data into training and testing dataset with a ratio as 4:1.

Thereafter, the model parameters were to be decided which again were chosen after a comparative study of those.

Several methods have been implemented to improve the performance of the algorithm.

- The value of k:
I have used the nearest neighbour approach to predict the rating of a movie by a user. The number of nearest neighbours is to be optimized to get the best result. I have collected the results for each iteration for different values of 'k' (number of neighbours).
- Distance metric:
Different distance metrics are applied to find the nearest neighbours.
 - Euclidean distance:
This is the most widely used distance metric. Euclidean distance is pretty straight forward and easy to calculate.
The Euclidean distance between two vectors X and Y is calculated as,

$$d(x, y) = \sqrt{\sum_i^n (x_i - y_i)^2}$$
 - Normalized Euclidean Distance:
It's always a good idea to normalize the rating for a user. For example, suppose one user tend to give a good rating to movies. So a rating of 6 out of 10 is very

poor. And the other user might be quite reserved when it comes to rate movies. So 6 could be average for him. So it's always better to normalize the ratings so that scale comes to equal.

➤ **Manhattan Distance:**

This is another distance metric when the absolute magnitude of the distance is desired.

The Manhattan distance between two vectors X and Y is calculated as,

$$d = \sum_{i=1}^n |x_i - y_i|$$

➤ **Minkowski Distance:**

This is the generalized version of both Manhattan and Euclidean distance. I have implemented this distance with the degree of 3.

The Minkowski distance between two vectors X and Y is calculated as,

$$d(\mathbf{x}, \mathbf{y}) = \left(\sum_{i=0}^{n-1} |x_i - y_i|^p \right)^{1/p}$$

➤ **Pearson Correlation Distance:**

This is the most useful distance metric when it comes to collaborative filtering. It takes care of the variance of the ratings for each user.

The Pearson correlation distance between two vectors X and Y is calculated as,

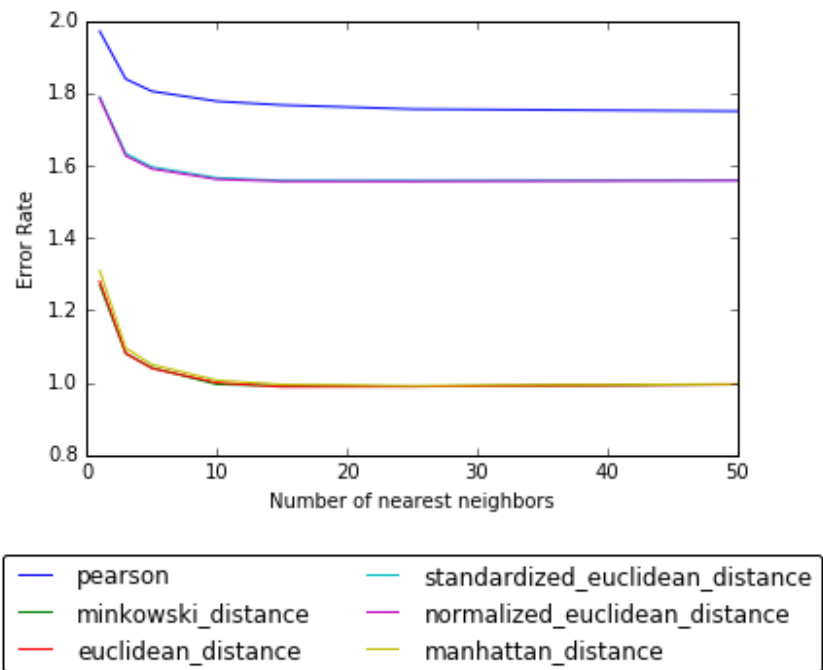
$$r(X, Y) = \frac{\frac{1}{n} \sum_i x_i y_i - \mu_X \mu_Y}{\sigma_X \sigma_Y}$$

where μ_X and μ_Y are the means of X and Y respectively, and σ_X and σ_Y are the standard deviations of X and Y. The numerator of the equation is called the covariance of X and Y.

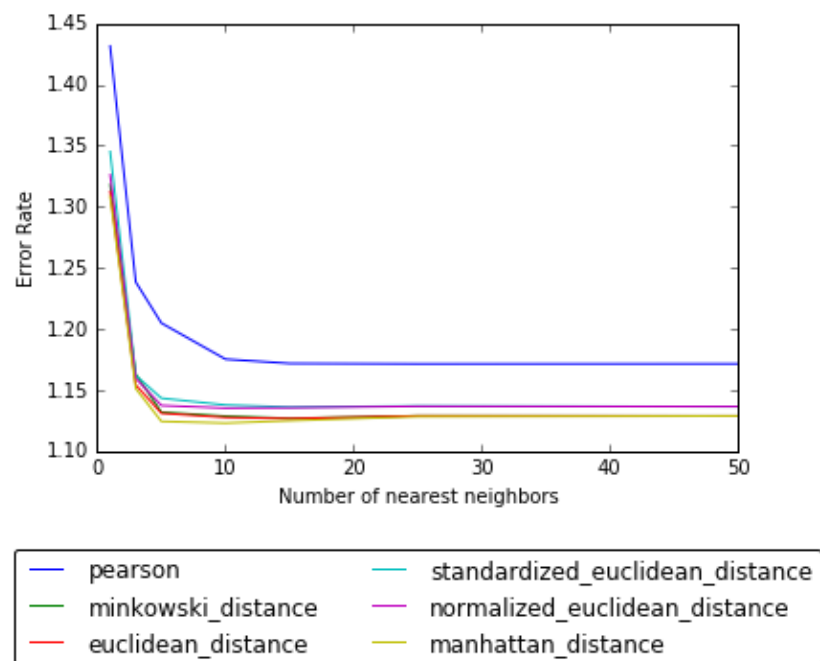
- I have implemented both the user based filtering and movies filtering to predict the rating and checked the accuracy. Also, the average of the previous 2 ratings are tested for accuracy.

IV. Results

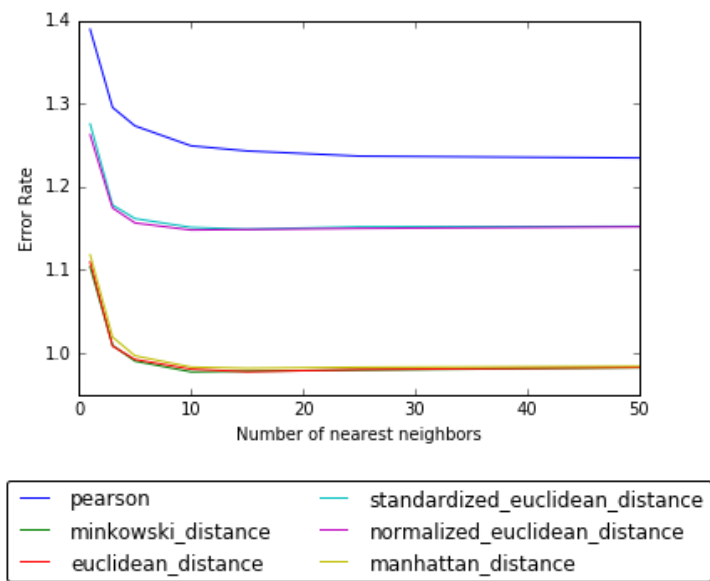
I have evaluated the recommendation engine with different parameters which are explained in the refinement section. Now let's check the comparative results.



Error Rate vs number of nearest neighbors for predicted ratings of user based collaborative filtering



Error Rate vs number of nearest neighbors for predicted ratings of movie based collaborative filtering



Error Rate vs number of nearest neighbors for average ratings of user based and movie based collaborative filtering

The performance of the recommendation engine is very good overall. Let us dig more into the parameters.

Number of nearest neighbours

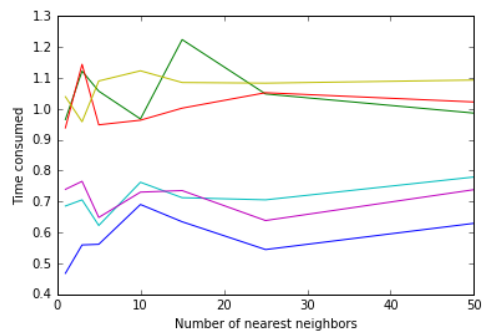
All three figures show a sharp decline in the error rate from $k=1$ to $k=3$. It further slowly decrease up till $k=10$ and then freezes. So we can conclude that $k=10$ is the optimal value of the number of nearest neighbours.

Distance metrics

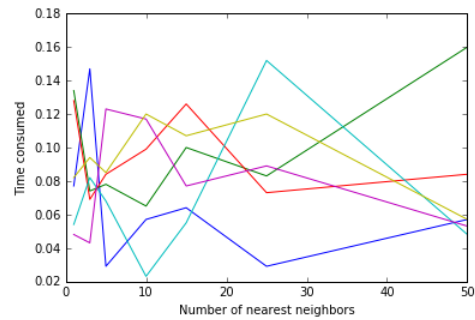
The above 3 figure clearly shows that the engine performs best when the distance metric is either Manhattan or Euclidean. These show significant difference in performance when used for user based filtering or the average of used based and movie based.

Pearson Correlation metric performs the worst as a distance metric which surprises me because this is one of the most widely used distance metric when it comes to collaborative filtering.

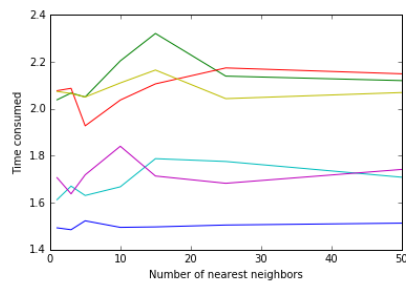
Now let's check the speed of the engine for different parameters. The below 3 figures clearly shows that the engines runs fastest when used with Pearson correlation distance. Manhattan, Euclidean and Minkowski make it slow. From the performance of the algorithm, we can say that there is trade-off between speed and performance when Manhattan, Euclidean or Pearson correlation are used. But the engine is fast enough; so we can choose Manhattan or Euclidean for better performance.



Time consumed vs number of nearest neighbors
for predicted ratings of user based collaborative filtering



Time consumed vs number of nearest neighbors
for predicted ratings of movie based collaborative filtering



Time consumed vs number of nearest neighbors
for average ratings of user based and movie based collaborative filtering

Justification

I had set the benchmark of RMSE for this algorithm as 1. The lowest RMSE I could achieve was 0.9819. Also, I discussed that a small difference in RMSE on the scale of 5 is actually pretty large. So I have achieved an RMSE 0.02 less than the target value which is quite good.

V. Conclusion

From the results section, we can see that both the speed and accuracy of the engine is quite good. The engine performs the best when used with Manhattan or Euclidean distance as the distance metric and with the number of nearest neighbours to be considered as 10. Also, we can see that it's wise to use the average of the ratings calculated from the movie based filtering and user based filtering.

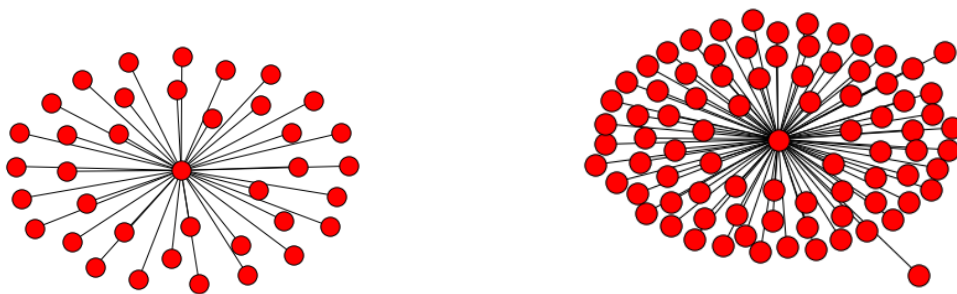
Reflection

It was a fun project to work on. The objective of the project was pretty straight forward and widely used in many applications. That makes it even more interesting. There were challenges as well for this kind of problem. There is no feature and instance dataset which is the most standard form of data in machine learning. Also I just had the rating data which made the approaches limited to collaborative filtering. But the size of the dataset is large enough to make the model quite general.

Free-Form Visualization

There is an interesting visualization I want to share. The users are not well connected when it comes to the common films they rate. I have built a network graph where the nodes are the users and the edges signify that there is at least one common film the connecting users have rated.

Let me show two networks; first one for the first 1000 ratings which contains only 38 users. Second one contain 10000 ratings which contain 84 users.



The above network clearly show that the users are sparsely connected. So while trying to find nearest neighbours, we don't have much competition. One way to deal with this could be calculating the rating based on an approach which is a combination of both user based and movie based collaborative filtering.

Improvement

I have used collaborative filtering to predict the rating of a movie by a user. There is another popular method used for the same purpose which is content based filtering. One of the major area of improvement for this engine is to implement the prediction based on content based filtering. Doing a comparative study between these two kinds of filtering for recommendation system would be a nice approach. The features for each movie for my dataset wasn't quite strong. It would be naïve to apply content based filtering for this dataset. But with enough information collected for the movies, we can perform a content based filtering for the same user and movie pairs.

Another improvement can be to implement the approach explained in the last section.