# Improvements to Collaborative Filtering Algorithms

by

Anuja Gokhale

_____

May 1999

APPROVED:

_____

Major Advisor

Prof. Mark Claypool

_____

Head of Department

Prof. Micha Hofri

# Abstract

The explosive growth of mailing lists, Web sites and Usenet news has caused information overload. It is no longer feasible to search through all the sources of information available in order to find those that are of interest to an individual user.

Collaborative filtering systems recommend items based upon opinions of people with similar tastes. Collaborative filtering overcomes some difficulties faced by traditional information filtering by eliminating the need for computers to understand the content of the items. Further, collaborative filtering can also recommend articles that are not similar in content to items rated in the past as long as like-minded users have rated the items. Unfortunately, collaborative filtering is not effective when there are too few users that have rated an item or for users that do not have a strong history or correlation with other users.

Content-based systems use content to filter or recommend items. These perform well when users know and specify topics in which they are interested. Recommendations for a user are based solely on a profile built by analyzing the content of the items which that user has rated in the past. Content-based filters face problems of over-specialization. When the system can only recommend items scoring highly against a user's profile, the user is restricted to seeing items similar to those that has already seen. Also, it is often difficult for content-based filters to understand the meaning of text or even the actual content of complex items.

We combine the strengths of content-based filtering techniques with collaborative filtering to provide more accurate recommendations. We use thresholds to improve the accuracy of traditional filtering algorithms, and design and implement a way to apply content-based filtering to an online newspaper. We compare our improved algorithms to current algorithms using both offline and online experiments and show that these result in more effective filters that can help manage the massive amount of information that is surrounding us today.

## 0.1 Acknowledgements

# Contents

# List of Figures

vii

# Chapter 1

# Introduction

Recent years have seen the explosive growth in the amount of information. The amount of information available through books, movies, news, television programming, advertisements, and in particular on-line sources, such as email, Usenet News, Web documents, is staggering. An individual cannot hope to absorb even a tiny fraction of today's information, and, to make things worse, more information is added daily.

The World Wide Web started in 1991 and grew to over 10,000 Web sites by December 1994 [RN94], and to 650,000 Web sites by January 1997 [Gra97]. Together these services comprise a minimum of 100 million documents. Furthermore, studies in 1995 have shown that the Web is filled with transient information. In 1995, there were an estimated 20 million pages on the Web and each page was online for only an average of 75 days. Other recent surveys have estimated the number of Web users in the U.S. as of May 1996 at 27,067,000 and the number of Web pages as of April 1996 at approximately 320 million [Ins96]. Furthermore, the number of pages is doubling every year. Using the average Web page size of 10 kilobytes (including graphics) brings the current size of the Web to 3 terabytes (or million megabytes).

Usenet News is also growing exponentially. Estimates show that in May 1996 there

to let the system give a stronger prediction even if the number of people who have read the article and have a correlation with the user is low. This is possible if the system bases the prediction on the content of the article also. We shall henceforth refer to this issue as the *Weak Rating Problem*.

Furthermore, current collaborative filtering systems do not use the content of the article as a factor in determining its relevance. Predictions are entirely based on correlation between other users. If a user wants articles on "collaborative filtering" to be given a high prediction, they cannot specify that articles with particular content (in this case "collaborative filtering") be always given a higher rating. This means that a user cannot setup specific profiles so that articles with certain keywords would be given higher rating predictions by the system. She is also unable to specify just how much she wants this information content to extend a rating prediction. For example, she may want items marked as at least average by other users and containing the words "collaborative filtering" to be given a higher prediction. Therefore, we also need to integrate current traditional content-based filtering techniques with collaborative filtering to provide a user with a more accurate prediction under all circumstances. We shall henceforth refer to this issue as *Content Problem*.

Lastly, tuning is also required for cases where current techniques provide predictions outside the range of valid ratings. Current collaborative filtering algorithms can face situations where the computed prediction falls outside the range of ratings that a user can give to an article. In such cases there is no clear definition of how to interpret these predictions that are above or below the valid range. For example, if a prediction of 7 is mapped to 5 (5 being the highest possible on a scale of 1 to 5), it is unclear how a user should interpret a prediction of 5 that was not re-mapped. This reduces the confidence that a user can place in a prediction. We shall henceforth refer to this issue as the *Normalization Problem*.

## 1.1 Contributions

We implement modifications to a novel weighted average algorithm to solve the problems outlined above. We improve accuracy in predictions arising due to the *History Problem* by enforcing thresholds on the number of articles seen in common by any two users. Ratings of users who have seen articles less than the threshold with the user in question are not considered for the computation of the prediction. We also impose similar thresholds on the correlation between two users so that only users who have very similar/dissimilar opinions influence the prediction. We analyze our modifications to the basic algorithm on data sets extracted from the *Jester* data set [RAJ97]. In addition, we analyze the characteristics of collaborative filtering data sets.

We derive techniques for content filtering of the news articles and for the integration of the collaborative and content filtering predictions. In order to analyze our approach to the content problem, we built a filtering system that combines collaborative and content based filtering on newspaper articles from a local online newspaper (*Telegram and Gazette*). We ran this system with a test bed of 8 users who rated newspaper articles for two months and collected and analyzed performance data.

In summary, the main contributions of our work are:

- We design and develop a more accurate collaborative filtering algorithm by implementing correlation and history thresholds on a basic algorithm.

- We show that our algorithm is also more consistently accurate.

- We show that the implementation of both history and user thresholds is most effective when these thresholds are user specific.

- We show that the integration of pure collaborative filtering with content based filtering improves the accuracy of the prediction especially in the startup phase

of a collaborative filtering system.

- We also show that the integration of collaborative filtering techniques with content-based filtering allows the collaborative filtering system to compute the predictions for new users or for other users who haven't rated enough articles to have a correlation with any other user.

- We show that the absolute size of the data set and the total number of ratings in the data set is not a good indicator of the accuracy of the collaborative filtering system. The accuracy of the predictions depends more on the number of rated items in common, because even there are either the total number of items in the system or the absolute number of items seen and rated by any user.

## 1.2  Outline

In Chapter 2, we look at related work in information filtering in general and collaborative filtering specifically. In Chapter 3, we discuss the various improvements we propose to the base algorithm and our experiments to evaluate our improvements. In Chapter 4, we discuss our approach to integrating collaborative filtering techniques with content-based filtering techniques and our experiments to evaluate our approach. Lastly, we present our conclusions in Chapter 6 and recommend areas for future work in Chapter 7.

# Chapter 2

# Related Work

The general problem of information overload have received considerable attention in research literature. Research in the general area of solving information overload problem can be broadly categorized into *information filtering techniques* and *information retrieval*. We shall use the term *information filtering* generically to refer both to finding desired information (filtering in) and eliminating that which is undesirable (filtering out).

## 2.1 Information Filtering Techniques

Malone et al. describe three categories of filtering techniques, cognitive, social, and economic, based on the information source the technique draws on in order to predict a user's reaction to an article [MGT+87]. The three categories provide a useful road map to other literature on filtering techniques. In the recent past, work has also been done on systems that use a combination of any two or all three of the above categories. We shall refer to these as *hybrid techniques* and shall discuss some of the work done in this area.

## 2.1.1 Cognitive Filtering Techniques

Cognitive, or content-based, filtering techniques select documents based on the text in them. For example, the kill files and citing result feature provided by Usenet news clients perform content filtering. More sophisticated techniques might also filter out articles from people who previously co-authored papers with an objectionable person. Even the division of Usenet news into newsgroups is a primitive example, since a reader restricts her attention to those articles with a particular text string in their "Newsgroups" field. Strings could also be combined with the boolean operators AND, OR, and NOT.

Alternatively, the profile of what to filter in or filter out could consist of weight vectors, with the weights expressing the relative importance of each of a set of terms [DDF+90, FD92, FS91]. In standard "keyword-matching" vector systems [S75], textual databases are represented by a "word-by-document" matrix whose entries represent the frequency of occurrence of a word in a document. The similarity between documents is computed as the inner product or cosine of the corresponding two columns of the word-by-document matrix. The words are considered to be pairwise independent.

Latent Semantic Indexing (LSI) [DDF+90, FDD+88] does not consider the words to be pairwise independent. In LSI, the associations among terms and documents are calculated and exploited in retrieval. A description of terms, documents, and user queries based on the underlying latent semantic structure is maintained. LSI can thus retrieve relevant documents even if they do not have any words in common with the query.

Some content filtering techniques update user profiles automatically based on feedback about whether the user likes the articles that the current profile selects. Information retrieval researchers refer to this process as relevance feedback [S88]. It has been shown that user inputs about concepts related to those mentioned in

an initial query, together with their relative importance, can significantly improve retrieval effectiveness [ ]. Relevance feedback can be improved if users select features from the texts of relevant documents [ ], instead of limiting them to selecting concepts from lists of terms selected automatically from relevant documents by the system (the system extracts these concepts by applying natural language processing techniques to the descriptions of interests that the user provided). The techniques for updating profiles can draw on Bayesian probability [ ], genetic algorithms [ ], or other machine learning techniques.

## 2.1.2 Social Filtering Techniques

Social filtering techniques select articles based on relationships between people and on their subjective judgments. A moderated newsgroup employs a primitive form of social filtering, choosing articles for all potential readers based on evaluations by a single person, the moderator. Collaborative filtering, based on the subjective evaluations of other readers, is an even more promising form of social filtering. Human readers do not share computers' difficulties with synonymy and context when judging the relevance of text. Moreover, items being filtered need not be amenable to parsing by a computer. People can judge texts on other dimensions such as quality, authoritativeness, or respectfulness.

The Tapestry system makes more sophisticated use of subjective evaluations [ ]. In Tapestry, many people can post evaluations, not just a single moderator, and readers can choose which evaluators to pay attention to. Moreover, filters can combine content-based criteria and subjective evaluations. For example, a reader could request recipes containing the words "hazardous substance" that users have evaluated and where the evaluation contains the word "good". Tapestry, though, does not include aggregate predictions.

The subjective evaluations used in collaborative filtering may be implicit rather

Rating Services have tools for rating and filtering Internet content. Some "customized newspapers" developed by using cognitive filtering techniques are Ad Hoc News, Scan News, CRAYON etc. An ongoing MCI at BYU is building a system to develop an on-line newspaper (Target) using social filtering techniques [MCL96].

## 2.1.3 Economic Filtering Techniques

Economic filtering techniques select articles based on the costs and benefits of producing and reading them. For example, Malone argue that mass mailings have a low production cost per addressee and should therefore be given lower priority [MGT87]. Applying this idea to Usenet news, a news client might filter out articles that had been cross-posted to several newsgroups. More complicated schemes could provide payments (in real money or reputation points) to readers to consider articles and payments to producers based on how much the readers liked the articles.

Stodolsky has proposed a scheme that combines social and economic filtering techniques [Sto90]. He proposes on-line publications where the publication decision ultimately rests with the author. During a preliminary publication period, other readers may post ratings of the article. The author may then withdraw the article, to avoid the cost to her reputation of publishing an article that is disliked.

## 2.1.4 Hybrid Techniques

The Ash system is a hybrid of the cognitive and social filtering techniques discussed above [BS97]. Ash maintains user profiles based on content analysis and directly compares the profiles to determine similar users for collaborative recommendation. Users "receive" items both when they score highly against their own profile or when they are highly rated by a user with a similar profile.

The Kbps system is similar to [BS97] except that during a similarity assessment between users, the system selects profiles of users with the highest correlation with an

individual user [CN92]. The Ringo system uses *mean-squared difference* (a measure of the variance of articles) and the *Pearson-R measure* (a measure of the inherent strength of the relationship between two sets of values) to determine similarity.

Another System [BHC98] also combines social and cognitive filtering techniques by defining additional features about the articles (for e.g. some features of a movie could be the the actors, directors, writers etc. of the movie). These features along with the proportion of a retrieved set of documents that are relevant (called *precision*) and the proportion of all relevant documents retrieved (called *recall*) are then used to categorize the articles by degree of likability for users. Correlations with actual users are not used in this system.

[SM'99] use filtering agents (filterbots) that act like normal users in a collaborative filtering system. These filterbots return ratings on articles based on certain syntactic information. This system used filterbots that generated ratings depending on the length of the article, the occurrence of spellings in the article etc.

The tech report by Microsoft Research compares the various collaborative filtering techniques and performs an empirical analysis on the same [BHK].

## 2.2 Information Retrieval

Conventional information retrieval (IR) [SM83] is very closely related to information filtering in that they both have the same goal of retrieving information relevant to what a user wants while minimizing the amount of irrelevant information retrieved [SM83, vR79]. Selective Dissemination of Information (SDI), one of the original information retrieval systems, is similar to most information filtering applications [IS78, LK70]. SDI was designed as an automatic way of keeping scientists informed of new documents published in their area of specialization. SDI maintained keyword based profiles of users and used these profiles to match the keywords against new

articles to predict which of the articles would be most relevant to the scientist's interests. Furthermore, research done in the field of evaluation of IR techniques [SM83, vR79] can also be applied to information filtering systems. A number of measures of evaluating IR techniques have been developed with the best known being precision and recall described above. These measures can also adequately evaluate the effectiveness of most information filtering techniques. Further, Belkin and Croft identify the primary differences between information filtering and retrieval [BC92]. This will help researchers in the area of information filtering to benefit from research in IR by "customizing" the IR techniques for information filtering while keeping in mind the differences between them. These differences mainly arise because user preferences in information filtering typically represent long term interests while queries in IR represent a short term interest that can be satisfied by performing the retrieval. Also, information filtering is typically applied to streams of incoming data while in IR, changes to the information contents do not occur often and retrieval is not limited to new items in the information source. Finally, filtering involves the process of "removing" data from the stream while IR involves the process of "finding" information in that stream.

## 2.5 Summary

In pure cognitive filtering techniques, only a very limited analysis of the content can be performed. Only certain aspects of the text can be evaluated and other aspects like aesthetic quality of multimedia information, style of language of text, other computer non-assessable information like the network factors (e.g. loading time) are completely ignored. Also filtering of items that do not match a profile effectively inhibit the user from being able to see articles on new topics (outside her profile) thus narrowing the scope of articles she sees. On the other hand, social filtering techniques suffer from

# Chapter 3

# Collaborative Filtering Improvements

We develop an algorithm to compute predictions that are more accurate (closer to the rating the user would give the movie) than those given by standard algorithms. In this chapter we describe the basic algorithm, describe the general design and setup for experiments to evaluate our algorithms and explain our incremental improvements.

## 3.1 Basic Algorithm

Although there is an increasingly strong demand for collaborative filtering techniques, only a few different algorithms have been proposed in the literature thus far [RIS '94], [USE], [MCT '97], [Maltz, Russell, 1895], [UPS]. Furthermore, most of these algorithms are based on simple predictive techniques that use a measure of agreement between users in order to make predictions. Like these previous researchers, we consider a collaborative filtering algorithm that uses a weighted average to compute predictions. We choose these algorithms not only because they constitute a large portion of the algorithms (including those used in commercial products) but also

because they intuitively generate predictions for users based on the similarity between the interest profile of that user and those of other users. These algorithms compute the similarities (correlations) between user profiles or compute correlations between users by looking at their history of agreement over articles rated in common.

The basic algorithm defines correlation as a measure of the degree of like-mindedness between pairs of users. We define likability as a measure of how much the user will like/dislike the movie. It is the amount by which her rating would be above or below her mean. The mean, here, is an average over all the user's ratings. It gives an indication of how "stringently" or "leniently" the user rates items. The basic algorithm uses the Pearson correlation coefficient to make full use of ratings between users that have different rating systems by adding the likability to the average rating the user gives her items to predict a rating for that particular user. For example, user A may rate all items between one to three where one is bad, two is average and three good while user B may rate all items between three and five where three is bad, four average and a five good. The algorithm adds the likability to the average of user A (two in this example) to predict a rating for user A.

The general formula to compute the likability for an article for a user by the basic algorithm is:

$$likability = \frac{\sum_i (corr_i \times (ratings_i - average_i))}{\sum_i (corr_i)}$$

$corr_i$ is the correlation of user $i$ with the user for whom the prediction is being computed. $ratings_i$ represents the rating submitted by user $i$ for the article for which the prediction is being computed. $average_i$ is the average rating (the average of all the ratings for all articles given by the user) for user $i$. $n$ is the total number of users in the system that have some correlation with user $i$. These $n$ users are those whose ratings are used in the calculation of the prediction for user $i$.

Consider an example to demonstrate how the formula is used to compute the prediction for an article for a user, X. The following table gives the ratings for an

article by two users and their individual correlations with the user $X$. We assume that the mean rating for the user $X$ is also known and equals 5 in this example.

|  | User A | User B |
|---|---|---|
| Correlation | .8 | -0.2 |
| Rating Given | 5 | 8 |
| Mean | 6 | 5 |

Mean for user $X$, $average_X = 5$

The prediction will be computed as:

$$
\text{Likeliness} = \frac{(.8 \ast (rating_A - average_A)) + (-0.2 \ast (rating_B - average_B))}{.8 + 0.2}
$$

$$
= \frac{(.8 \ast -0.2 \ast )}{1}
$$

$$
= 0.???
$$

$$
\text{Prediction} = (average_X) + \text{Likeliness}
$$

$$
= 5 + 0.???
$$

$$
= 5.???
$$

$$(3.1)$$

Our basic algorithm considers ratings given by all the users for a particular article in order to calculate a prediction. This includes those users who have very low correlation between them. A Pearson's correlation coefficient in the range $-0.2 <$ correlation $< 0.2$ is considered to be a low correlation between users. Similarly $0.2 <$ correlation $< 0.5$ or $0.5 <$ correlation $< 0.8$ is considered a moderate correlation and $-1 <$ correlation $< -0.8$ and $0.8 <$ correlation $< 1$ is considered a high correlation between any two users. Statistically, there is very little certainty that the ratings of users with low correlation between them follow either a pattern of similarity or dissimilarity. This suggests that the rating of such users there not be

allowed to have a bearing on the calculation of the prediction. A contribution of this thesis is to improve the accuracy of the basic algorithm by implementing a threshold on the correlation such that only users with a correlation above the threshold will be able to affect the prediction calculated by the system.

We can also see that the basic algorithm can produce inaccurate results if the number of users whose ratings are considered are below a certain minimum. As an example, consider the case where two users A and B agree on most things except one topic. We would ideally not like to consider just user B's opinion to compute a prediction for user A for that particular topic (this is not handled by our basic algorithm). On the other hand, if there are a hundred users (including user B) who think like user A, then the rating of user B for that particular topic would not adversely affect the prediction as the rating of user B contributes a much smaller portion towards the computation of the prediction.

The history (number of items/articles rated in common between two users) should influence the prediction in some way. As an example, consider the scenario where user B has watched a hundred movies in common with user A and has always agreed with user A while user C has watched just one movie in common with user A but has agreed on that one movie. In this case we cannot really be sure about user A being in agreement with user C on the next movie. In fact the more movies user A watches in common with user C, and agree, the greater the faith we have in the opinion of user C. There should be some way to differentiate between each users and give the rating of user B more weight than that of user C. We improve the accuracy of the basic algorithm by implementing thresholds on the history in common between two users.

## 3.2 Experiments

In this section, we present experiments that evaluate our proposed improvements. We shall first describe the general experimental design. We focus on the technique we use to calculate correlations between users and to compute a prediction for a particular article for a user. We then describe the experimental setup for our improvements focusing on how we compute the accuracy of a prediction.

## 3.2.1 Design

The design for our experiments on the improvements to the basic algorithm will comprise of off-line experiments. This means that we do not perform these experiments on real users using a live system but instead perform these experiments on both data simulated by us and on data from previous collaborative filtering experiments. The simulated data consists of both random and pseudo-random data sets generated by us that simulate the ratings provided by users in a real system.

The off-line experiments mainly consist of experiments on the data from the *MediaMagic* collaborative filtering service. The *MediaMagic* service was part of a research project at the Systems Research Center of Digital Equipment Corporation. The service was available for 19 months from February 1996 to September 1997. During that time the database grew to a fairly large size, containing ratings from 7399 users for 1059 movies. User ratings were recorded on a numerical six point scale. The data set is publicly available and can be obtained from the Digital Equipment Corporation (now Compaq).

As a part of the off-line design we maintain separate files which store the ratings for different movies given by a user and the correlations between all pairs of users in the system. The ratings for the movies are extracted from the *MediaMagic* data set. These

---

Philip Morris Paul Wendigass

ratings are stored as a matrix of users with ratings in a text file and are later used to generate correlations between users. The correlations between users are generated as *Pearson's correlation coefficient*. The system then uses these ratings and correlations and calculates the prediction for every movie for every user. The predictions are calculated using the formulae described in our individual improvements. This is done by "pulling out" the rating for the movie in question for a user. This means that when the system calculates the prediction for that movie for the user, it "ignores" the rating value for that specific instance from the matrix, thus simulating a condition where the user has not watched the movie earlier and therefore has given no rating for it. The prediction is then compared with the actual rating for the movie in question for that user (the value that was "pulled out" from the matrix during the computation of the prediction) to give an indication of the accuracy of the prediction. We describe this process in detail in the next section covering the experimental setup. This process is carried out for every movie for every user in the data set (for every point in the generated matrix) to get the average accuracy of the prediction for the data set.

We conduct similar experiments on random and pseudo random data. For these experiments, we generate random/pseudo random data to represent the ratings submitted by users for movies. All the other steps are the same as those carried out for the off-line experiments on the EachMovie data.

## 3.2.2 Accuracy of Predictions

The accuracy of a prediction can be measured by determining the number and severity of errors. An error is the difference between the rating prediction supplied by the system and the rating given by the user. We test our algorithms by checking the variance and percentage of errors produced by our algorithm. An error is serious if the difference in the predicted rating and that given by the user after rating the item is very large while an error is trivial if the predicted rating value and the actual rating

given by the user twice but since there was level or that regarding the prediction to the nearest level would not lead to any recognizable error by the user. The trivial errors are therefore a function of both the user interface used by the system and the scale of the ratings. If the acceptable range of ratings is the list of positive numbers between one to ten then an error of 0.3 for example might not get noticed by the user as the prediction will get truncated to the nearest integer. Consider the case where the rating given by a user was 2. If the system computed the prediction as a 4.7 (so that the error is 2.7) but truncated that off to 2 before showing it to the user then the user would not even notice an error if the prediction was shown as a number. On the other hand, if the user interface used a small bar to indicate the prediction then the error might have been more noticeable. As a second example, consider the situation where the system allows all positive numbers between one and hundred as valid ratings. In this case an error of 0.3 by the system would be less noticeable even on the smallest user interface. Thus a 10 point scale to give ratings and predictions with all integers between 1 and 10 as valid ratings (or predictions).

A "good prediction" will have low deviation of the prediction from the actual rating. The metrics that we will use to measure the accuracy of our algorithms include the following:

- **Mean Absolute Error** $\overline{x}$: Also called the arithmetic mean, mean absolute error is defined as the sum of measures observed divided by the number of observations. The lower the mean absolute error, the better the algorithm. For our experiments, we compute the mean absolute error between the rating given by the user and the prediction computed by the system.

- **Standard Deviation**: The standard deviation is defined as the square root of the average squared deviation from the mean. This is a measure of how consistently accurate the algorithm is. For our experiments we compute the standard deviation between the mean error for the data set and the error computed for each

individual prediction.

## 3.3  Specific Improvements

In this section, we describe and test the individual improvements to our basic algorithm. We also analyze the results of these experiments.

### 3.3.1  Improvement: Correlation Threshold

The basic algorithm computes a prediction by using the ratings submitted by all the users of the system for the article in question. This includes the ratings submitted by users who have a low correlation with the user in question. This can lead to inaccuracies in the prediction if the number of users with a low correlation increases or is more than the number of users with a high correlation with the user in question. The positive effect of the ratings of like-minded users on the prediction will be negated (or overshadowed) by the ratings of users with low correlation with the user in question.

As an example, consider the case where user A has a high correlation (equal to 0.9) with user B but a low correlation (equal to 0.1) with users C, D and E. If user B gives a rating of 4 (on a scale of 1 to 10) for an article and users C, D and E give a rating of 5, 6 and 7 respectively then the prediction computed by taking all the users into consideration will be 5.7 which is closer to 6 than to 4. Intuition suggests that since user A generally agrees with user B than with the other users, hence the prediction for user A should reflect the rating given by user B more than the ratings given by the other users.

Another potential problem with using the ratings of all the users is that there is an increase in the computation time which is not justified by an equivalent increase in accuracy. Considering the ratings of all the users (irrespective of their correlation with the user in question) increases the computation time exponentially with an increase

in the number of users. At the same time, this increase in computation time is not justified by a corresponding increase in the accuracy of the prediction. In fact, as we can see from the example above, the accuracy can potentially drop!

To address this problem, we implement thresholds on the correlations with users that are to be considered. For example, if we had only taken the ratings of users above a certain threshold (say 0.8) then we would have based our prediction that was closer to user $i$ than to the other users. We consider only those users that have a high correlation with the user in question thus reducing the margin of error. The advantage with this technique is that users with lower correlation do not unduly affect the rating prediction and the prediction is based only on those users who show either a high degree of agreement or disagreement.

The specific implementation of this improvement is as follows:

- Find the number of users who have a correlation greater than the threshold with the user in question.

- For all such users, apply the basic algorithm described in section 2.

We carried out off-line experiments on data extracted from the EachMovie database [EM97]. Data sets were created by extracting the ratings submitted by a subsection of the users and various movies. In order to obtain a broad range of correlations between users. These ratings were then stored as matrices with each element $(i, j)$ in the matrix (dataset) representing the rating submitted by a user $i$ for a particular movie $j$. We generated the prediction for every user for every movie. This was done by ignoring the element (rating) in the matrix for that particular user and movie. We then considered the ratings submitted by other users for the movie in question. We note that for this improvement we only considered users who have a correlation above the threshold with the user in question. We compared the prediction computed as above with the rating (extracted from the EachMovie data) that the user gave for
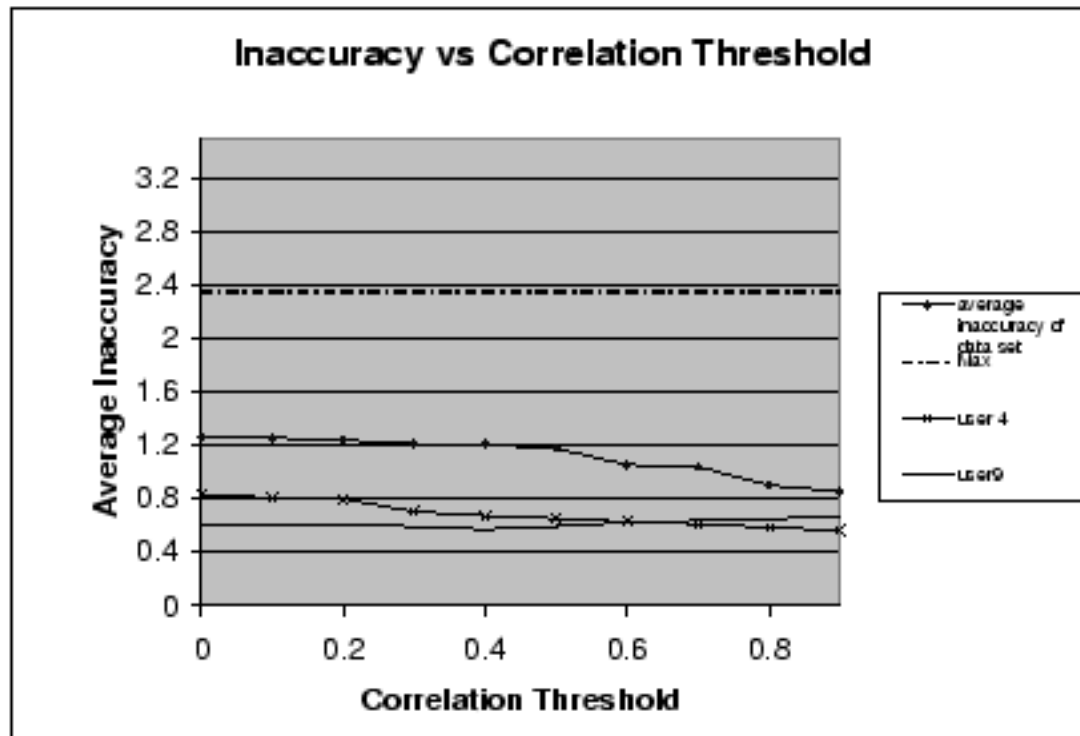
Figure 3.1: *Plot showing the relationship between the average inaccuracy and the correlation threshold for data from the EachMovie Data. This dataset consists of the ratings of 58 users and 125 movies. Average inaccuracy is an average of the differences between the rating and the prediction for every pair of users and movies in the data set.*

the movie. We repeated this same procedure for every user and movie (for every data element in the matrix data set). We computed the average absolute inaccuracy (absolute difference between the generated rating and its computed prediction) for the whole data set. This is the average inaccuracy of the data set for that particular value of the correlation threshold. We then repeated the same procedure for different values of the correlation threshold.

Figure 3.1 illustrates the results from our experiments on a data set consisting of ratings submitted by 58 users for 125 movies. The ratings in the data set are extracted from the EachMovie data. The ratings submitted were on a scale of 1 to

Figure 3.2: *Plot showing the relationship between the average inaccuracy and the correlation threshold for a random data set consisting of 10 users and 100 movies. Average inaccuracy is an average of the differences between the rating and the prediction for every pair of users and movies in the data set.*
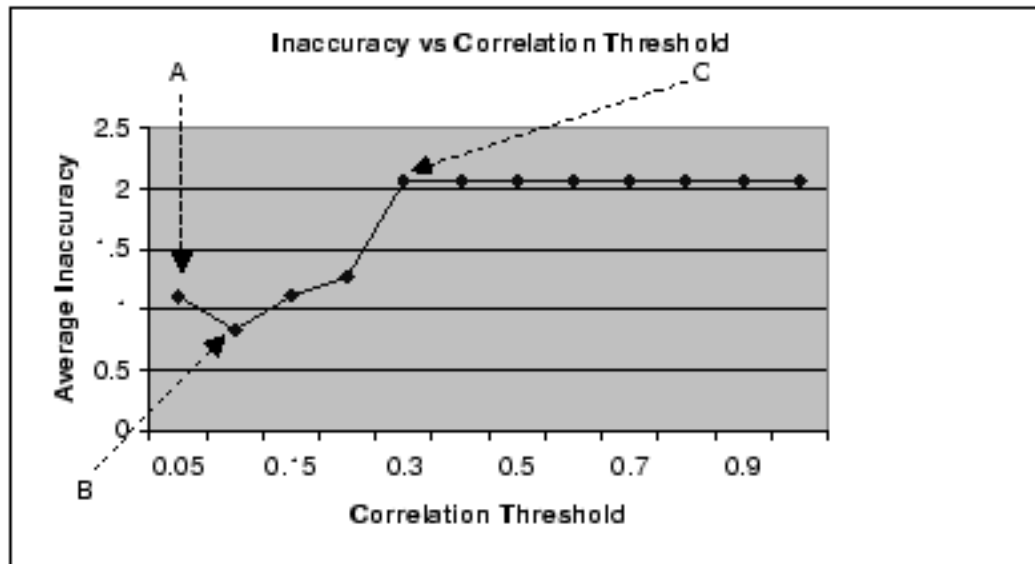
Figure 3.3: *Plot showing the relationship between the average inaccuracy and the correlation threshold for a random data set consisting of 10 users and 500 movies. Average inaccuracy is an average of the differences between the rating and the prediction for every pair of user and movie in the data set.*

the data set for the corresponding correlation threshold value.  The average inaccuracy was calculated by taking the average of inaccuracies in all the predictions computed by the system for every pair of users and ratings.

Point A on the curve shows the average inaccuracy when there is no correlation threshold applied (correlation threshold is 0.0).  Point B on the graph shows the average inaccuracy when the correlation threshold is 0.1.  Point C shows a rise in the inaccuracy when the correlation threshold is increased to 0.3.

Figure 3.3 also illustrates the relationship between the average inaccuracy of the data set and the correlation threshold.  This graph illustrates the results of the same experiment (with correlation thresholds) on a random data set with the same number of users but a greater number of movies (five times the number of movies as in the first graph).

We can see from the above graphs that the use of random data to test for collaborative filtering techniques is not very useful.  The inherent randomness of the data

generated does not lend itself to computations of correlations or the predictions made using the correlations. We make the following general observations:

- It is difficult to model real life data representing the opinions of people on various topics using random data sets. We normally expect people to have similar opinions on different movies/articles on the same topics. This behavior is difficult to reproduce in random data.

- It is also very difficult to "randomly" generate two sets of numbers such that the correlations between them remain almost the same or to generate them such that the next pair of numbers in each set models the correlation of the two sets up to that time.

- Finally, two sets of data generated randomly have almost no correlation, so that if these sets are used to represent ratings submitted by users then the data set would consist of "users" such that all the users have a very low correlation with each other.

We shall hereafter perform all our off-line tests only on data from the EachMovie data set [XX97].

## 3.3.2 Improvement: History Threshold

Typical collaborative filtering algorithms also use the ratings submitted by users who have a low history with the user in question. This means that the algorithms also consider the ratings of users who have rated a low number of items in common with the user for whom the prediction is being calculated. This can lead to inaccuracies in the predictions computed by the system if the number of users with a lower history with the user in question increases or is more than the number of users with a high history with the user in question. There is no statistical basis for a prediction if the

Figure 3.4: *Average inaccuracy and the history in common between any two users. The number of users is kept constant at 125. Average inaccuracy is an average of the differences between the rating and the prediction for every pair of users and movies in the data set.*
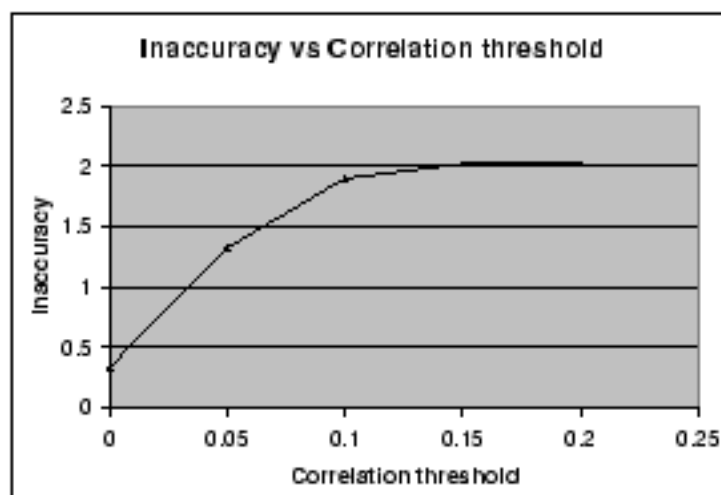
relationship of the inaccuracy of the prediction and the number of users in the data set for data extracted from the EachMovie database [EM97].

Figure 3.4 illustrates the results from our experiments on a data set consisting of ratings submitted by fifty eight users for six hundred movies. The ratings in the data set are extracted from the EachMovie data. The ratings submitted were on a scale of 1 to ten. The line labeled *Max* shows the average of the ($rating - mean$) of the data set. This is computed by taking the average of the inaccuracies in the prediction if the prediction returned was the same as the user's mean rating value. We plot this as this is the prediction returned for the user if no user in the data set has a history in common higher than threshold with the user in question. Line *average inaccuracy of data set* shows the average inaccuracy for the data set for the corresponding history threshold value. The average inaccuracy was calculated by taking the average of

Figure 3.5: *Plot showing the relationship between the average inaccuracy, the standard deviation of the average inaccuracy and the history in common between any two users in the data set for a data set consisting of ratings extracted from the EachMovie database. Average inaccuracy is an average of the differences between the rating and the prediction for every pair of user and movie in the data set.*

the inaccuracies in the predictions. We then plot the average inaccuracy of the data set and also plot graphs which are at 1 standard deviation distance from the average. (For normally distributed data approximately 60 percent of the total predictions would lie in the range within 1 standard deviation of the mean.

The graph shown below illustrates the relationship of the inaccuracy as well as the standard deviation of the inaccuracy of the prediction and the history in common between the users in the data set for data extracted from the EachMovie database.

We observe in figure 3.5 that the standard deviation is really high for a history in common of five movies and then rapidly drops with an increase in the history in common between users. This supports our earlier conclusion that an increase in the

history in common between any two users will lead to predictions that are not just more accurate but also more consistently accurate.

### 3.3.3 Improvement: Normalization

A potential problem with the basic algorithm is that the rating prediction computed by the system for a particular article for a user could be outside the valid range for the rating. This means that if the ratings given by the user are on a scale of 1 to ten, then the prediction could be either greater than 10 or less than 1 for some user. This can potentially lead to problems of interpretation of a prediction. For example, if the system consistently returns a prediction of 12 (while the allowed range of valid ratings is between 1 to ten) it will be forced to represent this prediction of 12 as a 10 as that is the maximum possible in the range. This leads to questions about the confidence a user can place in a prediction computed by the collaborative filtering system as the has no guarantees that if a prediction of 12 should be taken as a 10 then a prediction of 13 for example should still be taken as a 10 (since this is in the valid range and the system does not have to "bound" it off as in the previous case) and not as a 9.

We try and reduce the effects of this problem by normalizing the likability so that the final rating prediction is always between the minimum and maximum rating values. This allows us to convert the prediction value calculated into an accurate measure of the likability of a movie.

The specific improvement is as follows:

- Calculate the likability for user a using the formula in the basic algorithm:

$$likability_a = \frac{\sum_j (correlation\ of\ user\ j) \times ((rating_j - avrat_j)\ of\ j)}{\sum_j (correlation\ of\ user\ j)}$$

- $Rating_a = (avrat_a) + likability_a$

- let Max rating allowed $= q_{max}$ and minimum rating allowed $= q_{min}$

If (likelihood > 0) Normalise (maximum likelihood - likelihood ) to (q_max - mean)

If (likelihood < 0) Normalise (likelihood - minimum likelihood ) to (mean - q_min)

- Final rating prediction = mean + normalised likelihood value

We carried out various experiments using both random and pseudo-random data sets and the data from the EachMovie database. We also carried out various tests and computed the predictions on all combinations of movies and users for varying data sets. These predictions were computed for data sets of varying size and shape (ratio between number of users and movies weighted by users in a data set). We found, however, that though such a situation (of obtaining predictions above or below the valid range of ratings) can be theoretically possible, it never happened in all our tests. We conclude that this "improvement" need not be implemented with the intent of improving the quality of the prediction.

# Chapter 4

# Content-Collaborative Integration

In this chapter, we describe the motivation for an integration of content and collaborative filtering techniques. We also describe the content filtering algorithm and its various phases. We implement content-based filtering using a keyword matching technique. In the following sub-sections we shall talk about the techniques we use for keyword generation and matching and their application to the user profile in our content-based filter.

## 4.1 Motivation

The results from figure 3.3 suggest that collaborative filtering by itself cannot always guarantee a good prediction. On the contrary, the inaccuracy can increase if the number of people who have a correlation with the user in question is very low. In particular, figures 3.1, 3.2 and 3.3 suggest that correlation thresholds calculated for a data set produce more accurate predictions as long as the number of users which have a correlation above this threshold with the user in question is above a certain limit.

We performed experiments to observe the effect of the number of users on the average inaccuracy of a data set. We have seen earlier in figures 3.2 and 3.3 that the

46

Figure 4.1: *Plot showing the relationship between the average inaccuracy and the number of users in the data set for a data set consisting of ratings extracted from the EachMovie database. The number of articles is kept constant at 125. Average inaccuracy is an average of the differences between the rating and the prediction for every pair of users and movies in the data set.*

to be able to compute a correlation between them. In the beginning stages of the collaborative filtering system, the number of articles rated in common between any two pair of users is very low (Collaborative filtering systems are very sparsely populated). This can lead to situations where the system may not be able to compute a prediction at all or will compute the prediction based on the ratings of very few users making it highly prone to inaccuracies.

Also even if two users have a high correlation between them, they may not agree on particular topics. A user may want to see items on certain topics irrespective of the opinions of other users.

All the above observations suggest that pure collaborative filtering techniques are not sufficient and there is a need to make them more accurate. We suggest using content based filtering techniques in conjunction with pure collaborative filtering techniques to improve them.

In the following sections, we shall describe the content filtering algorithm and the integration of content and collaborative filtering techniques proposed by us.

## 4.2 Algorithm

### 4.2.1 Profile Format

We implement the content based algorithm using a keyword matching technique that relies on the significant words of the article and the user specified profile. We shall, therefore, first briefly describe the format of the user's profile required for the calculation of the content based prediction.

The user's profile is setup as follows. Each profile is divided into sections that the information source is divided into. These sections correspond ideally to newsgroups if the prediction is for a Usenet News system, to the various sections that a newspaper is divided into or the various categories that web sites are divided into

It would thus be more beneficial if the articles liked by the user in the past can also be used by the system to give a higher rating for articles that are similar (in our case, containing similar keywords) to those which the user gave a higher rating. This becomes more significant when we consider the fact that some users may not be able to give explicit keywords to describe their choices.

This implies that the system should consist of three parts and then combine them at the end. The three parts we refer to are:

- Implement a matching function on the explicit profile and the article keywords. We shall henceforth refer to this score as $M_1$.

- Implement a matching function on the user's list of implicit keywords and the article keywords. We shall henceforth refer to this score as $M_2$.

- Compute a score depending on whether the user wants to see all the articles in the section under which article in question appears. We shall henceforth refer to this score as $M_3$.

Until conclusively proven otherwise, we shall treat them all equally (we shall assign equal weights to them during the computation of the content-based score).

## 4.2.2  Key-Word Generation

As mentioned earlier, we base our algorithm for content-based filtering on the technique of keyword matching to determine how close the article matches the user's keywords. Since the aim of this research is just to explore the effect of integrating content-based filtering techniques with pure collaborative filtering and to design a way to implement the same, we feel it is sufficient to implement a simple content-based filtering algorithm like keyword matching. We base our algorithm on the technique proposed by Salton which states that the frequency of occurrence of words in an article

Figure 4.2: *Plot showing the relationship of the frequency of a word in an article and its rank order determining its importance in the article.*

furnishes a useful measurement of word significance [Luh58]. Let $f$ be the occurrence of various words in a text and $r$ their rank order, that is, the order of their frequency of occurrence. A plot relating $f$ and $r$ yields a curve similar to the hyperbolic in figure 4.2. This is in fact a curve demonstrating Zipf's Law [Zip49] which states that the product of the frequency of words and the rank order is approximately constant. Zipf verified his law on American Newspaper English. Luhn used this as a hypothesis to enable him to specify two cut-offs, an upper and a lower (see figure 2.1), thus excluding non-significant words. The words exceeding the upper cut-off are considered to be common and those below the lower cut-off rare, and therefore not contributing significantly to the content of the article.

Edmundson and Wyllys generalized Luhn's work by normalizing his measurements

suffix is removed from *looked* but not from the word *again*. To avoid erroneously removing suffixes, context rules are devised so that a suffix will be removed only if the the context is correct:

- The length of remaining stem exceeds a given number; the default in such a system is usually 2.

- The stem-ending satisfies a certain condition, e.g. does not end with a.

Many words, which are equivalent in the above sense, map to one morphological form by removing their suffixes. Others, unfortunately, though they are equivalent, do not (for example, *absorbing* and *absorption*). Since there is no easy way of making these fine distinctions, we put up with a certain proportion of errors and accept that they will not significantly degrade effectiveness. Experiments have shown that the error rate tends to be the order of 5 percent [Fra??]. It has also been shown that using a slightly different approach to stemming [Lov??] also produces errors of the same order of magnitude.

Once we have generated keywords for an article, we need to match sets of article keywords with those in the user profile and return a numeric indication of the closeness of the match of the article for the user. We shall hereafter refer to it as the *Matching Function*.

## 4.2.4 Matching Function

We need the matching function to compute the degree of match between the article keywords and the keywords in the user's profile. Let $D$ be the set of keywords representing the document (extracted or explained above) and $Q$ be the set of keywords representing the user profile. We use the *Overlap Coefficient* to measure the degree of match. This coefficient, $M$, does not take into account the sizes of $D$ and $Q$. This is important as the number of article keywords could be much longer than the keywords

in the user's profile. Also the implicit keywords for a user grow with time and can be much longer in size than the users profile. In such cases it is important to use this technique as what we need to check is not how many of the words in the two sets match, but what percentage of the profile's words match. For example both users A and B should get the same score if all the words in their profile are found in the article. This should not change even if user A had specified only 5 words in the profile and user B had specified 100. This technique is a normalised version of the simple matching coefficient. The formula to calculate the overlap coefficient is given below:

$$ M = \frac{P \cap Q}{\min P , Q} $$

(4.3)

## 4.3  Content Based Filtering Algorithm

We now precisely describe the complete content based filtering algorithm. To calculate the rating $M$ for every article we maintain a list of possible suffixes (for example, ion, ual, etc.). We shall hereafter refer to this list as the suffix list.

- For every new article in the system:

  Retrieve the article text.

  Strip off the suffixes of all the words in the article. In case of two or more matching suffixes, strip off the longest suffix. What remains are what are called stems of words (e.g. foot out of football).

  Count the number of occurrences of each such stem.

  Calculate the order of frequency of occurrence of the stems.

  Formulate a list of keywords for the article by taking all the words that fall in between the upper and lower cutoff points.

- For every pair of articles and users:

  If the user has not already rated the article:

    * Apply the Matching function to get $M_1$ and $M_2$. Let $D$ be the set of keywords representing the document (one extracted above), $Q$ be the set of keywords representing the user profile (explicit profile) and $R$ be the set of keywords representing the user implicit profile. We apply the Matching to $D$ and $Q$ to calculate $M_1$ and to $D$ and $R$ to calculate $M_2$. (Refer to section 4.5.x for a description of the matching function).

    * Calculate $M_0$. If the section under which the article falls is selected by the user then $M_0$ equals one, else $M_0$ equals zero.

    * Calculate the section based score. The section based score can be computed by taking the weighted average of $M_0$, $M_1$ and $M_2$. We shall give a weight of one-third for each of $M_0$, $M_1$, $M_2$ as we consider all three to be equally important

$$M = \frac{(1/3) * M_1 + (1/3) * M_2 + (1/3) M_0}{\cdot}$$

    (4.x)

  Where the user rates an article:

    * If the rating submitted is high (above 3) then add the article keywords to the user's implicit profile. We do this as a rating of 3 or above on a scale of one to ten means that the user has really liked the article. As the article keywords "reflect" the article, we can use these keywords as a measure of which keywords the user has liked in the past. We can then use these to determine how much a user would like a new article.

## 4.4 Integration with Content Based Filtering

After we calculate the predictions computed by the content based filtering and collaborative filtering algorithms respectively, we need to integrate these scores and return one number that is an indication of how well the user would like the article. In this section, we shall briefly describe the technique we use to integrate the content based and collaborative filtering predictions into one aggregate prediction.

The content based predictions would be more accurate in cases where the number of users or the history between users is low. The content based predictions for particular users may also be more accurate in general if the the number of users who are similar to the user in question is low. So if user A and user B agreed on most issues except one and the predictions for user A for an article on that issue was based only on user B we can see that the prediction made then will not be accurate. On the other hand, if this prediction was based on the opinions of many like-minded users the margin of error would be much less. The content based score would also be more accurate on topics where the user reacts to an articles irrespective of the opinions of others.

Similarly there are cases where the collaborative filtering score should be relied on more. This happens if the user has not specified enough keywords in his profile. The collaborative filtering score would also be more reliable for articles which have similar but not the same words/topics described in the users profile. This is where the user's past history of agreement with other users may help as these other users may have the similar words in their profile and his agreement in the past with other users would indicate that the article in question should get a high rating.

Both the collaborative filtering and content based scores are important but the extent of their importance towards the aggregate score (or prediction) is very user-specific. We therefore propose a system where the aggregate score is a weighted average of the collaborative filtering and the content based scores. We also ensure

that the weights attached to the collaborative filtering and the content based score are different for different users. These weights may also change over time to reflect the changes in a user's tastes.

We integrate by giving an equal weight to both the collaborative filtering and the content based score for all the users. We then adjust these weights according to the rating the user gives to every article. This is done by comparing both the collaborative filtering and the content based score respectively to the rating given by the user for the article. If the collaborative filtering score is closer to the rating given by the user then we adjust the weights to give the collaborative filtering score a higher weight than the content based score for future articles and vice versa. The exact value by which the weights are adjusted depend on the history. The formula for calculating this adjustment is specified in the algorithm steps given below. This is an ongoing process making the system self-learning. The increase or decrease in the previous weights is also a function of the number of articles the user has already seen. This means that a difference in the respective scores and the rating given by the user would lead to a greater change in the weights if the user had seen just two articles than if he had seen two hundred articles before the article in question. Though such a system is slow learning, as it adapts to changes in a user's tastes slowly, it also ensures that the system does not change drastically by ratings to articles that are unusual. For example, in cases where a user wants to rely only on his profile for articles on basketball and the article in question happens to be on basketball. In such a case, even if the article has received a low collaborative filtering score we don't want to change the weights by giving the content based score a very high weight compared to the collaborative filtering as the collaborative filtering score may be more accurate in general and only be bad for articles on basketball.

We now concisely explain the complete algorithm to implement the above.

### Integration Algorithm

- Initialization:

  final prediction = 0.5 * collaborative filtering score + 0.5 * content based filtering Score

- Each time a user returns a rating for an article:

  Check the collaborative filtering score, content based score and the rating returned by the user for the article.

  Check to see if the rating returned by the user is closer to the content based score or the collaborative filtering score.

  If the rating returned by the user is closer to the content based score:

  - new weight for content based score = old weight of content based score + (0.1 /number of ratings seen by user)
  - if (number of articles rated by user < 10) then new weight for content based score = old weight of content based score + 0.1.
  - new weight for collaborative filtering score = old Weight of collaborative filtering score - (0.1 /number of ratings seen by user)
  - if (number of articles rated by user < 10) then new weight for collaborative filtering score = old weight of collaborative filtering score - 0.1.

  If the rating returned by user is closer to the collaborative filtering score:

  - new weight for collaborative filtering score = old weight of collaborative filtering score + (0.1 /number of ratings seen by user)
  - if (number of articles rated by user < 10) then new weight for collaborative filtering score = old weight of collaborative filtering score + (0.1 /10)

- new weight for content based score = old weight of content based score - (0.1/number of messages seen by user)
- if Number of articles rated by user > 100 then new weight for content based score = old weight of content based score = 0.1.

## 4.5   Experiments

We realize that to test a collaborative filtering algorithm that also integrates content based filtering technique, it is essential to be able to get both the ratings of real users and the content of the articles they rate. The EachMovie database has data only on the ratings of the movies seen by the users. There is no information on what movie the rating is for. Movies are represented only by article numbers and not names. This makes it unusable for our experiments on the integration of content and collaborative filtering experiments. For this purpose we have built a system that gives the predictions for newspaper articles for the users of an online newspaper [MSN98]. This system allows the users to rate these articles and collects the ratings given. The system allows users to maintain profiles (containing sections and keywords for sections) specifying the users interests (if any). The system then computes correlations between users and gives predictions (for articles not yet rated by the user) based on both the collaborative filtering base algorithm described in section 4... and the content based filtering algorithm discussed earlier. We will use this system to collect the ratings given by the user and the predictions computed by the system to test the accuracy of predictions.

In this section we first describe the general design of our system. We shall focus on the technique used to calculate correlations between users and to compute a prediction for a particular article for a user. We then describe the experimental setup for our content based improvement focusing on how we compute the "accuracy" (or

effectiveness) of a prediction.

## 4.5.1 Design

The online experiments use real ratings submitted by real users on a live system [MSK...]. This system is designed to develop a customized on-line newspaper for the Telegram and Gazette. This uses a group of 9 users as a test-bed and collects ratings from these real users as opposed to our off-line design. We shall now briefly discuss the design of the system.

The system can be broadly divided into two parts: the front end and the back end. The front end consists mainly of the graphical interface that allows users to login and setup user profiles required for the content based filtering. User profiles are divided into the various sections of the newspaper and users can choose to specify keywords that reflect their interests. The front end is also responsible for showing users the predictions for various articles and collecting the ratings given by users for the same. These ratings are then uploaded to a database.

The back end of the system mainly consists of the database (where the ratings and cases are stored) and the algorithms which compute both the correlations between users and the predictions for an article. Correlations between users are re-computed similar to the off-line design (i.e. the system uses the *Pearson correlation coefficient* to calculate the correlations) once every day. The back end is also responsible for converting the individual content based and collaborative filtering cases to an integrated overall case using the algorithm developed by us. This integrated prediction is then used as a measure of likability of the article. The prediction is later compared with the actual rating that the user gives for the article. This is done to get a measure of the accuracy of the prediction. The front end shows a sorted list of articles in the order of the predictions each time a user logs into the system.

Figure 4.3: *Plot showing the average inaccuracy in the collaborative filtering score with time.*

a user is reading and the user has not specified anything in his profile. In this case, the system can use neither the implicit keywords nor the users profile to compute the content based score.

We can thus claim that the integration of content based techniques with pure collaborative filtering techniques can ensure that the system can generate predictions for almost all articles (except sometimes the first article) for all the users.

We also perform analysis on the whole data set to see the relative inaccuracies of the pure collaborative filtering score, the content based filtering score and a combination of the two. The graph in figure 4.4 shows the inaccuracy of the individual collaborative filtering and content filtering scores as well as the integrated score (the prediction returned to the user by the system).

We can see from figure 4.4 that the collaborative filtering score is more inaccurate in the beginning. This as we mentioned earlier is because of the the fact that the

Figure 4.4: *Plot showing the relationship between the inaccuracies in the content based filtering score, the collaborative filtering score and the integrated prediction. This dataset consists of the ratings of 18 users and all newspaper articles over 2 months. Percentage inaccuracy is an average of the percentage of inaccuracies in the prediction for every pair of users and movies in the data set.*

Figure 4.5: *Plot showing the relationship between the inaccuracies in the predictions computed using different techniques. This dataset consists of the ratings of 18 users and all newspaper articles over 2 months. Percentage inaccuracy is an average of the percentage of inaccuracies in the prediction for every pair of users and movies in the data set.*

- Our weighted average technique to integrate the prediction performs fairly well compared to the best a weighted average technique can perform.

The above results suggest that it is useful to integrate content based filtering with pure collaborative filtering. This is specially true as pure collaborative filtering techniques do not perform well in the initial stage and improve as users rate more items. In real life situations, this gets more difficult as most users stop using the system when they initially see inaccurate predictions or when the system cannot compute predictions for the particular user.

# Chapter 5

# Analysis of Results

- **Application of correlation thresholds leads to a corresponding decrease in the inaccuracy of the computation.**

From all our experimental results, we have reached the following conclusions:

There is a drop in the average inaccuracy of the data set with the increase in the correlation with users that are considered for the computation of the prediction, i.e. when the correlation between users is higher, the prediction computed is generally more accurate. We feel this is because of the fact that predictions made based on the opinions of people who have a low correlation with the user in question is more prone to error as these users do not have any established pattern of similarity/dissimilarity of opinions with the user in question. There is very little justification in considering the opinions of users with low correlation as their opinions on a particular item does not indicate how the the user in question might feel about the same item. Considering the opinions of such users, on the other hand, could potentially harm the accuracy of the prediction as their opinions could outweigh the opinions of users with a high correlation with the user in question.

The drop in inaccuracy is much more rapid in the initial increase in the correlation between users. When the correlation between users is above a certain value, the inaccuracy still drops with the increase in correlation between users, but this happens at a much slower rate. We believe that this is because of the fact that when the correlation between users is low, many users with a low correlation with the user in question could overshadow the opinions of the few users with a high correlation with the user in question. A prediction made in such a case, is highly prone to error. When the correlation between users increases, the margin of error in the prediction is lower as we are increasingly considering only those people who have similar/dissimilar to us. At this point, people who do not have patterns of similarity/dissimilarity are not considered at all.

- **Application of user thresholds leads to a corresponding decrease in the inaccuracy of the computation.**

We make the following observations from our experiments on user thresholds on this data:

There is a drop in the average inaccuracy of the data set with the increase in the number of users that are considered for the computation of the prediction. We feel this is because of the fact that predictions made based on the opinions of very few people can be more prone to error. Such predictions can also have a larger degree of inaccuracy if the few users it is based on do not agree with the user in question on that particular article even if they do agree with the user in question on most other topics.

The drop in inaccuracy is much more rapid in the initial increase in the number of users (in our experiments this is when the number of users is increased from the initial value of two users to a hundred users). When the

number of users is above a certain number the inaccuracy still drops with the increase in the number of users, but this happens at a much slower rate. We believe this is because of the fact that in the beginning, when the users in the system are very few, there are more cases of articles where the users who normally agree with the user in question may not agree with the same user on that particular topic. An example would be if user A and user B agree on most cases except on romantic comedies. Predictions for user A for articles on romantic comedies by taking just the rating of user B into account would lead to an inaccurate prediction. As the number of users in the system increases, such cases reduce and so there is a smaller drop in the decrease in the average inaccuracy of the system.

- **Application of history thresholds leads to a corresponding decrease in the inaccuracy of the computation.**

  We have conducted experiments on the data extracted from the EachMovie database. The data in the EachMovie database had been collected from real users using a live system that grew over time both in the number of users as well as the number of articles in the database.

  We make the following observations from our experiments on history thresholds on this data

  There is a drop in the average inaccuracy of the data set with the increase in the history in common between users that are considered for the computation of the prediction. i.e. when the number of articles seen and rated in common by users is larger the prediction computed is generally more accurate.

  We feel this is because of the fact that predictions made based on the opinion of people who have seen very few articles in common with the

the prediction made. This is the reason that the inaccuracy of the prediction drops at a much slower rate with any further increase in the history information between users.

- **Applying of history thresholds lead to predictions that are more consistently accurate**

We have seen that the standard deviation of the error in the prediction decreases with increasing history thresholds. As 68 percent of the total predictions lie within one standard deviation of the error, a decrease in the standard deviation means that 68 percent of the inaccuracies in the predictions lie within a small range of the average inaccuracy in the prediction. This implies that predictions computed are more consistent. As earlier results show that the inaccuracy decreases with the history thresholds, this implies that the algorithm gets more consistently accurate with increasing history thresholds.

- **Applying of normalization to the predictions is not very useful for real life data sets**

Though we have conducted experiments on varying sizes and shapes of data sets we haven't come across a single case where the prediction made by the system is outside the valid range of ratings allowed.

We therefore conclude, that such situations, though possible theoretically rarely occur. The implementation of normalization technique thus comes like a wasteful computation. We therefore have reached the conclusion that the implementation of normalization technique will not yield much benefit.

# Chapter 6

# Conclusions

The amount of information available to individuals is overwhelming. The Web, along with Usenet News, email, newspaper, books etc. contribute to an intimidatingly large information space. There is a clear demand for automated methods that filter as well as locate information with respect to users' individual preferences.

Collaborative filtering is a personalized technique to filter information. Collaborative filtering predicts the "likability" of an item based on the opinions of like-minded users in the system using their statistical correlations as the weights in a weighted average. These techniques consider the opinions of all users and therefore can be prone to inaccuracies in the prediction if the opinions of the few like-minded users get overshadowed by the opinions of many users with a low correlation. Users with low correlation do not have a pattern of similarity/dissimilarity of opinions with the user. Their opinions should ideally not be considered at all. Collaborative filtering also becomes inaccurate when the correlation between users is based on very few articles seen in common. Correlations based on a low number of articles in common are unreliable and may not correctly reflect the degree of similarity between two users. We reduce the collaborative filtering inaccuracy by implementing thresholds.

## 6.1 Correlations Threshold

The basic algorithm [GNOT92] computes the prediction by considering the ratings of all users that have a correlation with the user in question. We implement correlation thresholds such that we compute the prediction considering the ratings of only those users who have a correlation higher than the threshold.

We have shown that the application of correlation thresholds increases the accuracy of the prediction. There is, however, no correlation threshold that works best for all the users in the data set. Predictions are most accurate when correlation thresholds are specific to a user. This correlation threshold value depends on the number of users who have a correlation above this threshold value with the user in question.

A side benefit of this improvement is that correlation thresholds provide a reduction in the computation cost involved in the prediction as the number of users considered is smaller.

## 6.2 History Threshold

The basic algorithm [GNOT92] computes the prediction by considering the ratings of all users irrespective of the number of items seen and rated in common between the users. A problem with this technique is that users may not have seen enough messages in common to have the correlation between them accurately reflect their like-mindedness. For example, a history of only two messages in common and a correlation of 1 is probably not sufficient basis to say that the two users are likely to agree on the third message. We implement history thresholds such that we compute the prediction considering the ratings of only those users who have the number of items rated in common between them higher than the threshold.

Our experiments show that implementing history thresholds improve accuracy of the prediction. Although implementing a uniform history threshold increases the

accuracy of the predictions on average, the best value of the history threshold varies from user to user. Predictions computed for a user are most accurate when history thresholds are specific to a user. This threshold depends on the number of users who are considered for the prediction. The best threshold, therefore, depends on the number of users who have a history above the threshold with the user in question.

We have also shown that the application of history thresholds also leads to predictions that are also more consistently accurate. The standard deviation of error in the predictions decreases rapidly with increasing history thresholds making the error in the predictions more consistent. This observation can be used to get an indication of the strength of the prediction and can be returned to users to give them an indication of the confidence they can place in any particular prediction.

## 6.3    Integration with Content-based Filtering

Collaborative filtering provides inaccurate predictions to users that do not have many like-minded users. Also, collaborative filtering does not allow a user to specify that a particular content be given a higher prediction. Finally, with collaborative filtering alone, users who are the first to rate an item cannot get predictions for it. We integrate content based techniques with collaborative filtering to alleviate the effect of the above problems.

We devise an algorithm that implements keyword matching to perform content based filtering for news articles. Our algorithm uses the keywords specified in the user's profile and the articles keywords to compute a content based score. Our algorithm also uses a set of implicit keywords (extracted from articles liked by the user in the past) and the news categories (for example, world news, sports etc. ) that the user has selected to compute different content based scores. We have derived a technique to integrate these three scores (explicit keywords, implicit keywords and

the category) to compute an overall content based prediction for the user.

We use a weighted average that is adaptable to the relative accuracy of each score and user specific. The weights assigned to the collaborative and content based predictions are not only dynamically modified to reflect the user's choices but are also user specific. This is important as the collaborative filtering prediction might be more accurate for certain users while the content based filtering score might be more accurate for others, depending upon the correlation and history in common with other users as well as how well the user has set up her profile.

Our experiments show that the collaborative filtering score is less accurate in the startup phase of the collaborative filtering system. In extreme cases, the system is unable to compute a collaborative filtering score for a user. This can happen if the user is the first person to rate the article or if the system is unable to compute a correlation for the user with any other user (a user has to have rated at least two articles in common with another user for the system to be able to compute a correlation between them). We also see that in this startup phase the integration of the collaborative filtering score with a content based score leads to predictions that are more accurate. At the very least, it allows the system to compute some prediction based on the content based score if the system is unable to compute a pure collaborative filtering score. We have also shown that the integration of content based filtering technique allows the system to compute a score for new users who have not rated a single article. The only case when the system is unable to compute a prediction is when the prediction is being computed for a new user who has not specified anything in the profile and is reading the first article.

# Chapter 7

# Future Work

We implemented each of the improvements individually and show that they return a benefit over the basic algorithm. We have seen that correlation thresholds do not improve the accuracy of predictions unless the number of users is above a certain minimum. This, for example, suggests that the implementation of correlation thresholds in conjunction with user thresholds would be more beneficial. Though most of the individual improvements lead to an increase in the accuracy, we feel that implementing all the improvements would lead to more accurate predictions.

For our off-line experiments we have assumed that all articles belong to the same information space. That is, the correlation between users that we compute and consider for the prediction is based on all the articles that both the users have rated. Let us assume that two users have rated articles on different topics and have different degrees of correlation for each topic/subspace of the information space. A pair of users can have a high correlation on certain topics and a really low one on others. In this case, the system does not efficiently use the high correlation even if it is computing a prediction for an article in the subspace where the users have a high correlation. For example, consider the case where two users A and B always agree on the movies they see (have a high correlation for movies) but do not agree on most other topics (they

have a low overall correlation). In this case, though user A has a low correlation with B in general, user A would still like the system to give a higher weight to the rating of user B for all articles on movies. We do not implement non-linear correlation in our off-line experiments. We believe that computing separate correlations for different information categories for a pair of users will be more effective than computing one overall correlation between users.

Our experiments use the ratings given by a pair of users for all the articles/movies seen in common to compute a correlation between them. This includes those articles/movies that are universally liked (liked by most users in the data set) or disliked. It can be argued that such articles/movies are not really as useful in capturing the similarity between users as less universally liked/disliked ones. We can extract the actual similarity between two users' opinions if we compute the correlation based on just those articles/movies that are not universally liked or disliked. Such a technique might improve the accuracy with which the correlation captures the degree of similarity between people and thus improve the prediction.

In our content based filtering algorithm we perform a keyword matching of the user's profile and the article's keywords. We generate the article's keywords from it's text. Words falling between a lower and upper cutoff (refer section 4.2.2) are considered to be the keywords of the article. In our work, these cut-offs have been set to one-fourth and three-fourths of the rank order of the words. We feel that a study on the determination of these cut-offs and thus choosing the best cut-offs might improve the performance of the content based filtering algorithm thus improving the performance of the integrated algorithm.

For keyword matching in the content based algorithm we use the stems of words (the part of the word that remains after the word is stripped off it's suffix). Most words reduce to the same stem when the longest suffix is removed from them. For example, both "coats" and "coated" reduce to the stem "coat" when the longest suffixes in the

respective words, s and sk, are removed). Some words, on the other hand reduce to different stems even when they are essentially the same. For example, "subscription" and "subscribing" reduce to the stems "subscrip" and "subscrib" and will be treated as occurrence of different words. One way to reduce this problem would be to establish equivalent stem endings and then use these to identify similar words. In this case, "subscrib" and "subscrip" are equivalent stem endings. Implementing these techniques should lead to somewhat more accurate content based scores and therefore better predictions.

The accuracy of the keyword matching technique also depends on how accurately the keywords extracted from an item reflect its content. Removal of common words (such as articles) by pre-processing the content with a 'stop-list' before extracting keywords for the article might improve the accuracy of the keyword matching technique to compute a content based prediction.

We used the keyword matching technique on the user's profile and the article keywords to compute the content based score. One disadvantage of the keyword matching technique is that it depends on the exact word (the word "stamp" in our system) and does not exploit the semantic content of the keywords. This can lead to inaccuracies as people use different words to to describe the same concept. Latent semantic indexing (LSI) overcomes this problem of variability in human word choice by automatically organizing textual information into a semantic structure into a semantic structure more appropriate for information retrieval. Using LSI to compute the content based score might lead to more accurate content based score as LSI does not depend on exact word matching and will therefore perform better content analysis.

The complete content based algorithm in section ... computes three different scores to compute the content based score. These scores, $S_1$, $S_2$ and $S_3$ are based on the keyword matching between the user's profile and the explicit and implicit set of

keywords for the article. These scores are then integrated using a weighted average to compute the complete standard based score. The weights given to each of these scores by us are equal and are arbitrarily assigned. Future work in the correct assignment of weights to these scores will improve the standard based score and thus the integrated prediction.

We integrate the standard based and the collaborative filtering predictions using a weighted average. Both the predictions are assigned equal weights in the beginning. These weights are dynamically modified depending upon which prediction is closer to the actual rating submitted by the user. A future work in this direction is to investigate hill climbing technique as a means of adjusting the collaborative and standard-based weights.

One limitation with the Pearson's correlation coefficient we use is that it depicts a linear relationship. If a pair of users have a strong correlation but it is nonlinear then this may not be reflected as a strong correlation and used as such by the system. For example, consider two users A and B such that they have a low correlation when A likes an item but have a high correlation whenever A dislikes the item. This means that they dislike similar items but may not like similar items. The system should ideally be able to use this information so that it is treated as a "similar" user for A whenever A gives a low rating for an article. When calculating a prediction for an article for user A such that this article has been given a low rating by user B, the rating given by user B should be given a higher weight when the system computes the prediction.

Our collaborative filtering system cannot compute a collaborative filtering score for a new article that no users in the system have read. The new article has not been rated by any users and therefore the system does not have any ratings for the article making it impossible to compute the prediction for this article for any user. At least one user has to first rate an article before the system can calculate the collaborative

filtering score for any other user for that article. The system handles such situations by just returning a standard based score based on the user profile and the article's keywords. Such situations can be dealt with if the system also maintains correlations between articles. For example, if two articles get a similar rating by the same user then this means that the two articles are similar. Articles are accordingly be similar in standard if the standard based score returned for them are similar. Thus when a new article enters the system, users can be given a prediction for the article by checking which articles are similar to the new article and using the ratings given by the user for all such articles to calculate the prediction.

# Appendix A

## A.1 Sparsity Vs Inaccuracy with Time

The goal of collaborative filtering systems is to help people focus on reading documents that would be of interest to them. This requires either filtering documents that are not relevant or computing a score related to each document that represents how much the user would like the document. Our system takes the latter approach.

We have seen in our experiments on the MovieLens data that the accuracy of such a score for a document depends a lot on the number of people who have seen articles in common with the user for whom the score is computed. It also depends on the number of articles seen in common between the user in question and any other users whose rating is considered for the computation of the score. Both the above mentioned factors imply that collaborative systems work well when lots of users have rated lots of articles in common. The accuracy of the prediction (or score) may depend on the sparsity of the data set (we define sparsity as the total number of ratings in the system vs the product of the total number users and articles).

We conducted experiments to study the effects of the sparsity of the MovieLens data set over time and the corresponding accuracy of the data set.

We ran our experiments to observe the change in sparsity of the data set over time. Any data set grows over time as the number of articles coming into the system

Figure A.1: *This graph shows a possible distribution of ratings in a data set where some users have given ratings for all articles in the data set. Each point (x,y) represents the fact that a rating has been submitted by the user x for the article y.*

## Horizontal Distribution of Ratngs

Figure A.2: *This graph shows a possible distribution of ratings in a data set where some articles have been given ratings by all users in the data set. Each point (x,y) represents the fact that a rating has been submitted by the user x for the article y.*

Figure A.3: *This graph shows a possible distribution of ratings in a data set where some articles have been given ratings by all users in the data set and some users have given ratings for all the articles in the data set. Each point (x,y) represents the fact that a rating has been submitted by the user x for the article y.*
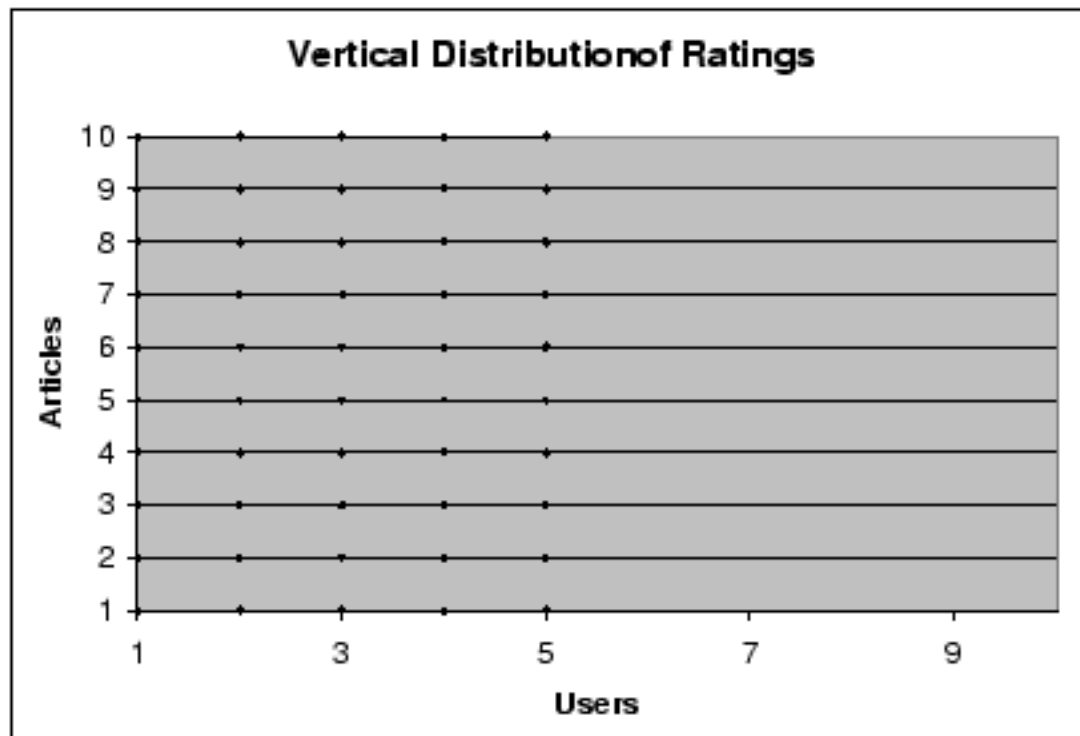
Figure A.4: This *graph shows a possible distribution of ratings in a data set where all the users have given ratings for some articles (not necessarily the same articles) in the data set. Each point (x,y) represents the fact that a rating has been submitted by the user x for the article y.*

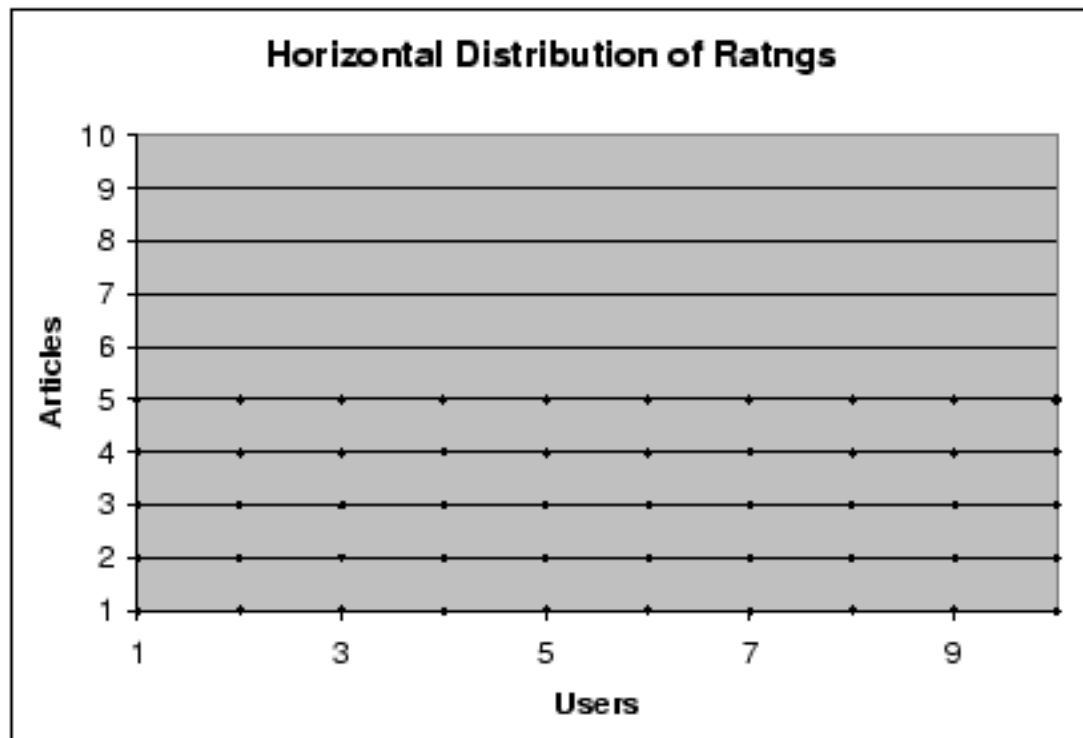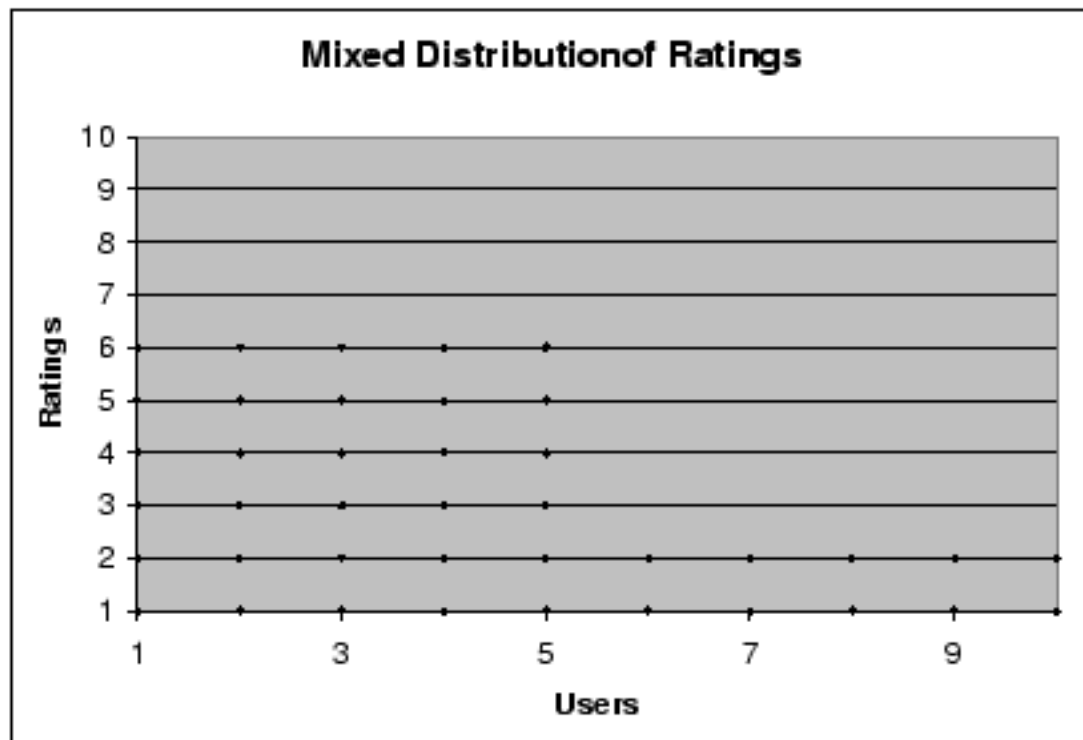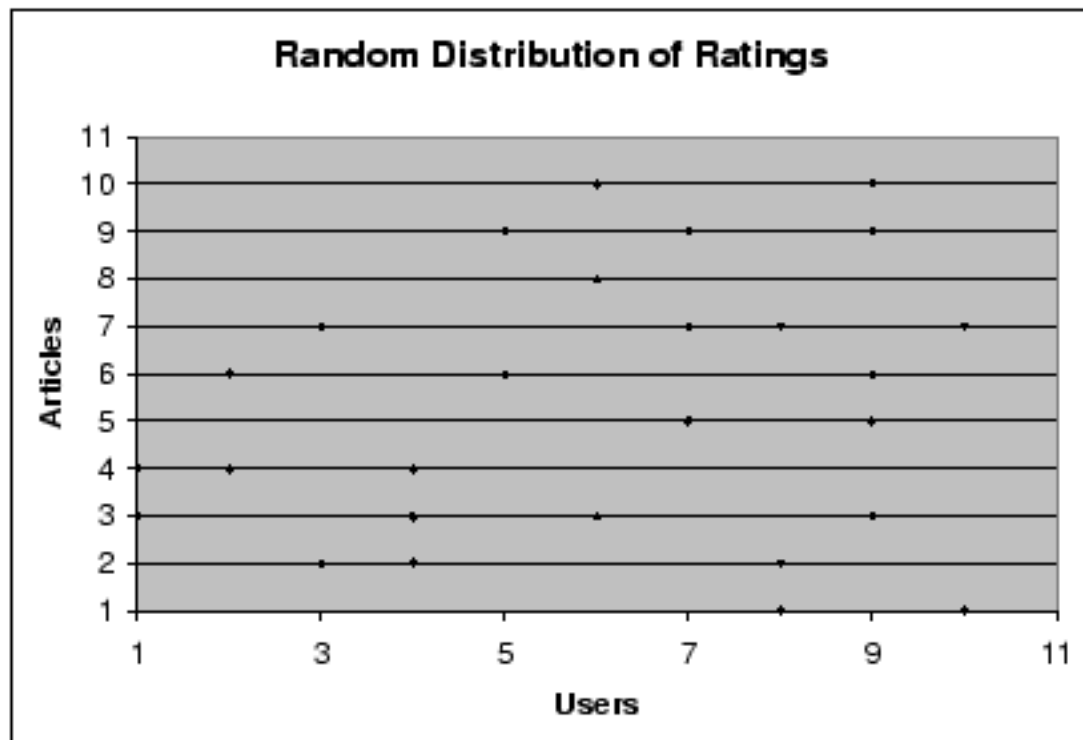Also some users may not have any articles in common with another user and so will have no correlation value with this other user.

We ran experiments to observe the effect of the sparsity of the data set on the average accuracy of the data set. These experiments were carried out to see if the accuracy of the data set increased with a decrease in the sparsity (higher ratio of actual number of ratings to total possible ratings in the data set) or was affected by the factors stated above.

Figure A.1 and A.2 show the sparsity and inaccuracy of the data set with time. We can see from figure A.1 and A.2 that both these values (sparsity and inaccuracy) of the data set are not a function of time. They neither increase or decrease steadily with time. Also the inaccuracy of the data set does not change inversely with the sparsity. We can thus safely say that the inaccuracy of the data set is not a function of the absolute sparsity of the data set. The data set does not get more accurate as more ratings are entered into the data set. Rather, the inaccuracy depends on the way these ratings are distributed. This means that a collaborative filtering system may not get more accurate with the addition of more users or articles.

This supports our hypothesis that the accuracy of the prediction is a function of the history in common between users and the number of users any particular user has a correlation with. This, in turn, depends upon the way the ratings are "distributed" in the data set.

## A.2 Effect of Shape of Data Set

For data sets where the number of users is large and the number of articles read in common is above a certain threshold, application of thresholds yields little benefit in the accuracy of the prediction. This is typically true for data sets which are "vertical" in shape and have very few (or none) of the users having a high correlation

Figure A.5: *This graph represents the average inaccuracy of the data set over time. The average inaccuracy of the data set is calculated by taking the average of the inaccuracies in the predictions for every user for every article.*
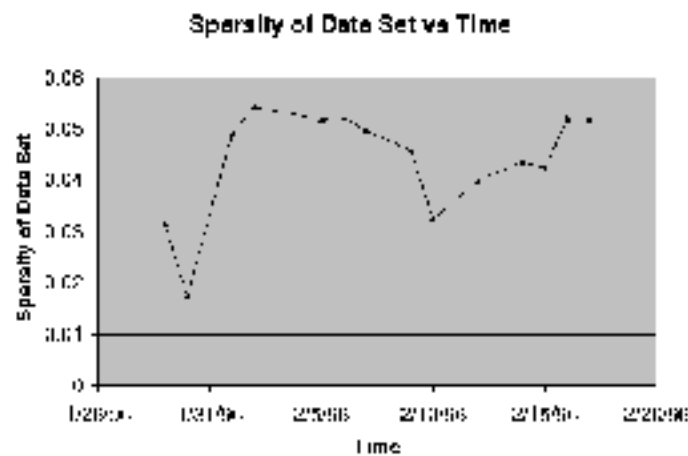


Figure A.6: *This graph represents the sparsity of the data set over time. The sparsity of the data set is calculated as the ratio of the total number of ratings in the data set and the total number of possible ratings in the data set.*

with another user. That is, for data sets where the ratio of users to articles and the number of articles read in common is big. This is so as such data sets typically have a low correlation among users as the users have seen a lot of articles in common and so the chance that any two users agree on almost all the articles (which would lead to a high correlation) is low. The application of thresholds though does not adversely affect the prediction (so long as the number of such users having a correlation above the threshold is reasonable) but leads to a benefit in the computation time.

We have found that the implementation of correlation thresholds yields to an increase in the accuracy for data sets that are horizontal in nature. That is, for data sets where the ratio of user to articles read in common by users is low. This is because users have a low to medium correlation with other users. Applying correlation threshold for such data sets excludes the users having a low correlation with the user in question. This leads to both a benefit in the accuracy of the prediction and the computation time. We also note that the application of correlation thresholds without the application of user thresholds can sometimes have adverse effects as the accuracy increase if the number of users above the threshold drops below a certain value.

# Appendix B

## B.1 Second Integration Algorithm

- At time t=0 final prediction = 0.5 * collaborative filtering score + 0.5 * content based filtering Score

- Each time a user returns a rating for an article.

    Check the collaborative filtering score, content based score and the rating returned by the user for the article.

    Calculate tempWeight1 = old Weight of content based score + (content based score-rating)/(collaborative filtering score-rating) (content based score-rating).

    Calculate tempWeight2 = old Weight of content based score + (collaborative filtering score-rating)/(collaborative filtering score-rating) (content based score-rating).

    new weight for content based score = Average of the tempWeight1 and all the distinct weights in the past.

    new weight for collaborative filtering score = Average of the tempWeight2 and all the distinct weights in the past.

# Bibliography

[And7.]    K. Andrews. The development of a fast condition algorithm for english. *Dissertation submitted for the Diploma in Computer Science*, 97..

[BC92]    N.J. Belkin and W.B. Croft. Information filtering and information retrieval: Two sides of the same coin? *Communications of the ACM*, 35(12):29–38, 1992.

[BHC98]    Chumki Basu, Haym Hirsh, and William Cohen. Recommendation as classification: Using social and content-based information in recommendation. *American Association for Artificial Intelligence*, page 7 – 728, 1998.

[BHK]    John S. Breese, David Heckerman, and Carl Kadie. Empirical Analysis of Predictive Algorithms for Collaborative Filtering. *Technical report*.

[BP98]    Daniel Billsus and Michael Pazzani. Learning Collaborative Information Filters. *Machine Learning: Proceedings of the Fifteenth International Conference*, 1998.

[BS97]    Marko Balabanovic and Yoav Shoham. Content-based Collaborative Recommendation. *Communications of the ACM*, 40(3), March 1997.

[C...] W. B. Croft and Das. Experiments with query acquisition and use in document retrieval systems. In *ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 349–368. ACM, 1993.

[DDF...] S. Deerwester, S. T. Dumais, G. W. Furnas, T. K. Landauer, and R. Harshman. Indexing by Latent Semantic Analysis. *Journal of the American Society for Information Science*, pages 391–407, 1990.

[DM97] DealMemo web site. http://www.DealMemo.com. Last Accessed, 1997.

[F...] G. W. Foltz and S. T. Dumais. Personalized Information Delivery: An Analysis of Information Filtering Methods. *Communications of the ACM*, 35(12):51–60, 1992.

[FDD...] G. W. Furnas, S. Deerwester, S. T. Dumais, T. K. Landauer, R. A. Harshman, L. A. Streeter, and K. E. Lochbaum. Information retrieval using a singular value decomposition model of latent semantic structure. In *Proceedings of the 11th Conference on Research and Development in IR*, pages 465–480. ACM/SIGIR, June 1988.

[Gray] Matthew Gray. Web Growth Summary. http://www.mit.edu:8001/people/mkgray/net/web-growth-summary.html, 1996.

[HW...] Henderson, H. R. and Wyllys R. A. Automatic abstracting and indexing: survey and recommendations. *Communications of the ACM*, 1958–59, 1958.

[HHPM...] W. C. Hill, J. D. Hollan, D. Wroblewski, and T. McCandless. Edit Wear and Read Wear. pages 3–9. CHI '92 Conference on Human Factors in Computing Systems, 1992.

[HTSM+95] W.C. Hill, L.D. Stead, D. Wichlewski, and T. McCandless. Using Collaborative Filtering to Weave an Information Tapestry. *Communications of the ACM*, 35(12):1–70, 1992.

[JK70] A.M. Hauptman and D.A. Kochala. State of the art of selective dissemination of information. *IEEE Trans. Eng. Writing Speech III*, 1:79–95, 1970.

[KMM+97] Joseph A. Konstan, Bradley N. Miller, David Maltz, Jonathan L. Herlocker, Lee R. Gordon, and John Riedl. GroupLens: Applying Collaborative Filtering to Usenet News. *Communications of the ACM*, 40(3):77–87, March 1997.

[Lev79] R.I. Levine. Silver Evaluation for clustering algorithms or clustering algorithms. *Journal of the American Society of Information Science*, 30(3):87–94, 1979.

[Luh58] H.P. Luhn. The automatic creation of literature abstracts. *IBM Journal of Research and Development*, 2:159–165, 1958.

[Mal94] D.A. Maltz. Distributing Information for Collaborative Filtering on Usenet Net News. MIT Department of EECS MS Thesis, May 1994.

[MGT+97] P.W. Maloney, K.J. Chard, J.A. Dubok, S.A. Jacket, and M.D. Cohen. Intelligent Information Sharing Systems. *Communications of the ACM*, 40(3):234–235, 1997.

[MSS+98] Tim Miranda, Matt Savino, Dmitry Yakar, and David Yuralkov. Collaborative Filtering. 1998.

[PR97] R.S. Pindyck and D.L. Rubinfeld. *Economic Models and Economic Forecasts*. McGraw-Hill, New York, 4th edition, 1997.

[Int99]     Internet Archive. http://www.archive.org. Last Accessed, 1999.

[Res94]     P.A. Resnick. GroupLens: A Collaborative Filter for Network News. MIT Department of EECS MS Thesis, February 1994.

[RIAA '97]  Peter Burman, Will Hill, Brian Amento, David McDonald, and Josh Gomez. Phoaks: A system for Sharing Recommendations. Communications of the ACM, 40(3), 1997.

[Rea99a]    Wire Trace.     www.opennowhat.com/individuals/99-01.htm.     Last Accessed, http://www.opennowhat.com/individuals/99-01.htm, June 1999.

[Rea99b]    Wire Trace. www.opennowhat.com/individuals/99-02.htm. Last Accessed, www.opennowhat.com/individuals/99-02.htm, May 1999.

[TW99]      Shardanand U and Maes P. Algorithms for automating 'word of mouth'. Conference on Human Factors in Computing Systems - CHI 1999, 1999.

[vR79]      C.J. van Rijsbergen. Information Retrieval. Butterworths, London, 2 edition, 1979.

[Zip49]     H.P. Zipf. Human Behavior and the Principle of Least Effort. Addison-Wesley, Cambridge, Massachusetts, 1949.