

Fuzzy Book Recommendation System

Evgeniya Evlogieva

Student @ Universiteit van Amsterdam
eevlogieva94@gmail.com

Philip Bouman

Student @ Universiteit van Amsterdam
10668667

Abstract—Nowadays, with the rapidly increasing number of e-books out there it is becoming easier to start reading a new book, it is only one click away. With the growing diversity, however, it is becoming harder to find a book that you will like. In this paper, we are presenting a book recommendation system that is using a Fuzzy Logic System (FLS). We will use a dataset of books and user ratings for these books. The FLS will take as input the genre membership of the books and the user ratings, and will give as output to what extent a user likes a certain genre. For measuring the performance of the system, we will conduct a survey amongst a group of about 20 students with different backgrounds, and ask them how would they grade our solution. We are planning to further extend the proposed solution by adding a FLS for matching users, and also by cluster the set of genres and introduce fuzzy membership to a certain cluster.

Index Terms—Fuzzy Logic, Recommender system, Book genre, Book recommendation

I. INTRODUCTION

Fuzzy Logic is one of the pillars of Computational Intelligence and it deals with uncertainties in real world problems. It finds application in a lot of domains such as control [1], decision making [2], [3], image processing [4], etc.

In our project we investigate the application of Fuzzy Logic in the domain of decision making, and, more precisely, in a book recommendation system. There have been several attempts to solve this problem [5], [6], however, none of them using Fuzzy Logic. The aim of recommendation systems is recommending items that the user would like. Whether a user would like something or not, and to what extent, is a matter of personal preference and uncertainty.

The purpose of the project is to recommend a list of books (and an indication of the extent to which the user will like a certain book), based on an input in the form of a list of book genres the user has rated.

The significance of the problem is not in its nature, but in the particular Fuzzy Logic approach, as it hasn't been applied yet to this particular area.

Our approach is to use an existing data set with user preferences and books. For the books for these datasets, we need to obtain their genres and to which extent a book belongs to a certain genre. This puts a certain limitation on the number of books we can use, because the 'scraping' of the genres is time-consuming for a rather big dataset.

Because of the above mentioned limitation we will use only a part of the data set as a training data (approximately 50 000 books), and a part of it as test data (the same amount).

The rest of the paper is organised as follows. Section 1A mentions our objectives and the approach for the evaluation

phase. Section 2 gives a brief outline of the current approaches for finding user preferences. Section 3 presents in detail our approach of the problem. In the different subsections we discuss the data that we use for testing, the design and the implementation of the FLS. Section 4 and 5 explain what experiments we held and what results we obtained. In Section 6 we give our personal critical opinion on the results. Finally, concluding remarks and future plans for extending the project are included in Section 7.

A. Objectives

The goal and our objective is to achieve at least 50% success rate. This means that on our test data, at least 50% of our "users" must actually have "liked" or "loved" any of the books our system recommends based on the part of their ratings we use as input.

Another measurement of success would be to have a survey and ask people to try the system. Considering the existing limitations, we will evaluate which approach will be more suitable for our case.

II. LITERATURE REVIEW

Understanding user preferences is an important part of recommendation systems. Most approaches in obtaining preferences rely on explicit feedback from users such as ratings or lists of interest. However this kind of feedback tends to be affected by user inconsistencies (*natural noise*). [7]

Another way to get information about user preferences is using implicit elicitation methods, which can automatically learn preferences from item features, users demographic information and past behavior such as web and purchase history. Different implicit methods can be distinguished: collaborative filtering, content-based, and hybrid approaches. [8]

Collaborative filtering methods approximate user preferences based on items rated by users, by finding similarities between users and recommend items that similar users liked but the target user had not come across yet. A downfall of the collaborative approach is that it does not take item features (genre, author) in to account and can be faced with a number of computational problems (scalability). Other limitations of collaborative filtering are that it cannot be applied to new items that have not been rated, and to unassigned users (*cold start/first rater problem, sparsity*). [9]

In a content-based approach user preferences are modelled (using machine learning) based on item features of user rated items. Automatically extracting item features can be hard

for certain items, since it requires careful feature selection, representation, and inference. Content-based methods also cannot predict users new behavior like having new interest (limited by users history). [8], [9]

A hybrid approach tries to combine both implicit feedback methods, using both item content and user rating behavior.

III. APPROACH

A. Data

There are two sources of data used. The first one is an existing dataset with data from bookcrossing.com [10] containing three different files: a list of user accounts, book information and user ratings for those books. The second source is a web crawl of Goodreads [11] consisting of a list of books and their genre's as classified by users from Goodreads.

The first dataset, created by combining the user account and user rating files from the first source, will consist of: *user-ID*, *user age*, and a list of (*book:rating*)-tuples as rated by this user. The books will be referenced by their International Standard Book Number (ISBN) and the associated rating will be on a scale of 1 to 10, or 0 if the rating is implicit. The original files contain information about 278858 user accounts and 1149779 rated books. The number of user accounts is reduced considerably since about two-thirds of the users did not rate any books and are therefore discarded from the final dataset.

(As of now the final dataset has not been constructed yet, so we can't give exact information about the size.)

An example of an entry (first dataset):

```
[user-ID, user age, (ISBN, rating), ... ]

[114, 57, ('0312953453', 7),
 ('0446608653', 9), ('0446612545', 9),
 ('0446612618', 8), ('0451208080', 8),
 ('0553584383', 9), ('0671027360', 10),
 ('0812575954', 5)]
```

The second dataset that will be used consists of genre information from Goodreads. By using the ISBN's from the book information data, a webcrawl was utilized to gather genre information about those books. As of now this set contains 33181 books and their genre's as classified by Goodreads users. The aim is to get genre information for approximately 50000 books. The genre information consists of a list of genre names (with at least one vote) and the number of times people rated it as such. Some pre-processing steps are applied, the most important one being the normalization of the number of votes. The final dataset will consist of an ISBN followed by a number of (*genre:grade*)-tuples.

An example of an entry (second dataset):

```
ISBN;[(genre, grade), ... ]

8445071408;[('fantasy', 1),
 ('classics', 0.48),
 ('fiction', 0.43),
```

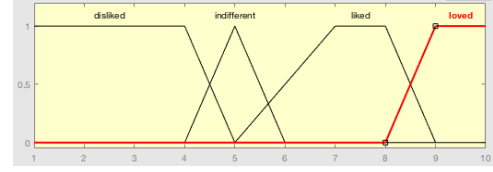


Fig. 1: Membership functions of user ratings.

```
("adventure", 0.1),
("science-fiction-fantasy", 0.08)]
```

The second dataset can be used as a lookup table for the entries in the first dataset. E.g.: find the appropriate genre information for a book rated by a user.

The output data of the system will consist of a list of ISBN recommendations in order of "best-match". The user will be shown this list along with accompanying book title and author.

B. Design

The Fuzzy Logic Inference System will have the following role: for each user, we will obtain a list of tuples of the type (*genre:grade*), describing how does the given user like a certain book genre. The recommendation part will happen later, and it will be based on similarities between users, calculated with an error function.

Now, about the input of the FLS. Our system works with input consisting of sets of book ratings for each user. Items are considered "loved", "liked", "indifferent" or "disliked" based on the given rating. This is where fuzzification is applied. On Figure 1 you see the membership function of the ratings to each of the fuzzy sets "loved", "liked", "indifferent" and "disliked". The range is from 1 to 10, because the ratings in the data set are in this interval.

For the "disliked" fuzzy set we have a trapezoidal Membership Function (MF), starting from 1, and ending at 5. This decision is taken, based on our reasoning, that every rating from 1 to 4 you give, you more or less disliked the book. The "indifferent" fuzzy set has a triangular MF, because the "truly" indifferent grade for a book you could give is a 5. The function is also symmetrical, because any value that is equally far from a 5 has to have equal membership degree. For a rating of 7 and 8, we can definitely say that someone liked the book. That's why the MF of the "liked" fuzzy set is trapezoidal with a core [7, 8]. It starts from 5 so that we have more overlap between it and the "indifferent" one, and it ends at 9. Lastly, the "loved" MF is also trapezoidal, with a core [9, 10].

Our second input is, for every book, a list of tuples (*genre:membership_to_the_genre*). The membership we calculated with a special normalization formula. Originally, in our data we had how many people 'shelved' a certain book as a certain genre. But when someone 'shelves', or grades a book, it does not give just one genre. That way, if we want to normalize the data by summing up all the people to 100%, we may have people, who were counted twice or three times. Therefore, after some tuning of the parameters, we chose the following solution: we will assume, that when a person is asked about

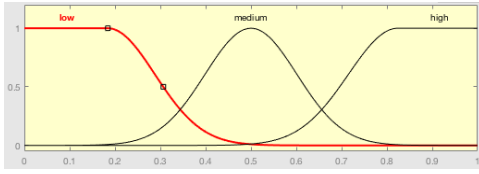


Fig. 2: Membership functions of genre membership.

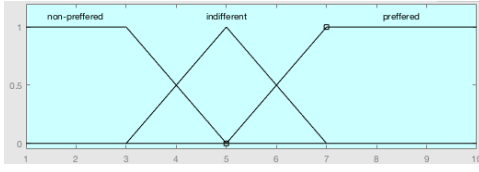


Fig. 3: Output - Membership functions of user preferences.

the genre of a book it will give, on average, three different genres. That's why, the normalization algorithm performed on our data will be the following: we sum up all the people that rated, or 'shelved' the particular book. Then, we divide that number by 3 to obtain the presumably real total number of people that rated the book. To obtain the real membership of a book to a genre, we divide the number of people who 'shelved' is as that genre, to the presumably total number of people. That is also beneficial, because in that way, you could have a book with a membership degree of one. In other words, the sum of the membership degrees of the different genres do not necessarily add up to 1. So having this membership of a book to a genre we fuzzify it to see how much it belongs to the fuzzy sets "low", "medium" and "high" (Figure 2).

For these three fuzzy sets we decided we wanted the membership functions to be smoother, so we chose the Gaussian shaped ones. The "low" and "high" are symmetrical to one another. The "medium" is with a core at the value 0.1. The range is between 0 and 1, because the normalized values that we have obtained are in this interval.

The next step in the FLS is the rules. For this step we decided to come up with the rules ourselves, based on general knowledge.

The output fuzzy sets are "preferred", "indifferent", and "non-preferred"(Figure 3). Much like the second input, the "preferred", and "non-preferred" MFs are symmetrical, and the "indifferent" is with a core at 5.

The crisp output of the FLS, is, for a certain user, a list of tuples (*genre:grade*), that presents how much this user likes the certain genre.

*Note: In the future, a process of tuning the parameters will be performed, so the above mentioned MFs might change in shape or range.

For the actual recommendation step we will calculate an error function (the absolute distance vector) between the target user and every other user that we have data for. The user with the lowest value of error will be the one we will get our recommendations from.

C. Implementation

For the implementation Python and MatLab will be used. The fuzzy inference part of the recommender system will be implemented in Matlab using the Fuzzy Toolbox. The final recommender system will be a Python script, taking user genre preferences as input and giving a list of recommendations as output.

A Matlab script will be used to setup and call the different parts of the fuzzy inference system within the Fuzzy Toolbox. With the Python **mlabwrap** package this Matlab process can then be controlled from Python to get the values for all the users and genre's.

To collect and preprocess the data a number of python scripts are used aswell. First **parser.py** is used to collect data from Goodreads [11]. Next **genres_memberships.py** is used to filter the appropriate data from the web crawl and normalize the genre votes. The **UserSet.py** script is used to combine the user account and user rating datasets from the Book-Crossing Dataset [10].

A prototype of the system was created by hand using the Matlab Fuzzy Toolbox, containing the input and output membership functions and the rulebase. (**Prototype.fis**)

All the code can be found on <https://github.com/phielp/Fuzzy>

IV. EXPERIMENTS (EXPERIMENT PLANS)

We plan to conduct an experiment with about 50 000 books and their genres. For these books, we will have a look at the user ratings. We will "run" our FLS on that data. Then, to begin with, we will try it out ourselves. We will go through a stage of tuning the system parameters. Also, for this stage, we might use some book recommendation sites to manually test a few scenarios. If we are happy with the results, we might set up a survey amongst some fellow students to collect their opinion on our performance.

Another experiment to measure the performance will be to get some of the other data (for now the plan is to use another 50 000 books), and based on some of the user ratings, to try to predict what the other user ratings are. This part of the evaluation experiment is still a bit fuzzy and it will be cleared out at a later stage.

V. RESULTS

VI. DISCUSSION

VII. CONCLUSION

APPENDIX A

SOME FORMULAS HERE

ACKNOWLEDGMENT

REFERENCES

- [1] Chuen-Chien Lee. Fuzzy logic in control systems: fuzzy logic controller. ii. *IEEE Transactions on systems, man, and cybernetics*, 20(2):419–435, 1990.
- [2] Hung-Tso Lin. Fuzzy application in service quality analysis: An empirical study. *Expert systems with Applications*, 37(1):517–526, 2010.

- [3] Mohammad Hossein Fazel Zarandi, Neda Mohammadhasan, and Susan Bastani. A fuzzy rule-based expert system for evaluating intellectual capital. *Advances in Fuzzy Systems*, 2012:7, 2012.
- [4] Gopala Sainarayanan, R Nagarajan, and Sazali Yaacob. Fuzzy image processing scheme for autonomous navigation of human blind. *Applied Soft Computing*, 7(1):257–264, 2007.
- [5] Zan Huang, Wingyan Chung, Thian-Huat Ong, and Hsinchun Chen. A graph-based recommender system for digital library. In *Proceedings of the 2nd ACM/IEEE-CS joint conference on Digital libraries*, pages 65–73. ACM, 2002.
- [6] Dong Kun. Research of personalized book recommender system of university library based on collaborative filter [j]. *New technology of library and information service*, 11:44–47, 2011.
- [7] Denis Parra-Santander and Xavier Amatriain. Walk the talk: Analyzing the relation between implicit and explicit feedback for preference elicitation. 2011.
- [8] Azene Zenebe, Lina Zhou, and Anthony F Norcio. User preferences discovery using fuzzy models. *Fuzzy Sets and Systems*, 161(23):3044–3063, 2010.
- [9] Anthony F. Norcio Azene Zenebe. Fuzzy modeling for item recommender systems or a fuzzy theoretic method for recommender systems. 2008.
- [10] <http://www2.informatik.uni-freiburg.de/~cziegler/BX/>. Book-crossing dataset, 2004.
- [11] <http://www.goodreads.com>.