

ECEN 765-600
Machine Learning with Networks

Academic Project Report

**Recommendation System based on Collaborative
Filtering with focus on Cold Start**

Submitted by:

Ankit Yadav(UIN 726006435)
ankityadav270796@tamu.edu

Project Guide: Dr. Xiaoning Qian



**ELECTRICAL & COMPUTER
ENGINEERING**
T E X A S A & M U N I V E R S I T Y

Department of Electrical and Computer Engineering
Texas A&M University
College Station, TX – 77843

Abstract

This work has been done as a part of academic project for subject **“Machine Learning with Networks”**. A lot of research has been performed on the recommendation system based on collaborative filtering, but cold start issue seems to have been ignored in most of the studies. Hence, I wish to dedicate this project to address the cold start problem. The purpose of this project is not to improvise any existing techniques but to identify if the user’s demographic or entropy0 scores can have any impact while addressing the cold start problem in collaborative filtering system.

First a basic recommendation system was implemented, it is being used as the base model. Any of the implemented model is compared to this model to see if there has been any improvement. Overall 4 models were implemented including the basic one.

Also, apart from the project implementation I have dedicated time to learn about the recent research work being done on “Cold Start problem in Recommendation System” and identify the areas where this project can be further taken to.

Introduction

Over the last decade, there has been a rapid growth in the internet users. The growth in internet users between 2000-2018 is more than 1000 percent and it's still growing. Companies like Amazon, Flipkart, eBay, Netflix and YouTube introduce new stuff every single day.

There is a high probability of the user being deluged with the content getting uploaded into these websites if the content was not personalized, and soon users would lose their interest. Here comes the recommendation system for the rescue. In terms of volume recommendation system is able to handle a big amount of data. Although there has been a good amount of research on this topic, but still this remains a topic of interest given the significance.

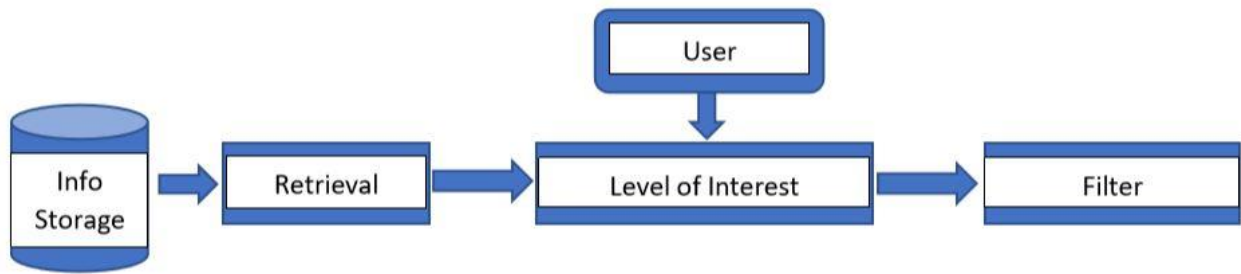
While looking for the similar projects, I noticed that most of them have not tried to consider the "Cold Start Problem" i.e. showing recommendation for a new user or new product that our system has no data available for. In this project, I expect to deal with that part of the problem. Apart from this, my plan is to do a literature survey on "Cold Approach Problem" in the recommendation system.

1.1 Recommendation System

Recommendation system suggests items to a user based on his/her past behaviors, their personal data, their physical location or based on some form of similarity to others. I seem to date back the presence of commendation system in my life since I had my first mobile phone in 2004, it was capable to show me favorite contacts based on who I talk to more or message.

Therefore, a recommendation system is a system that is capable of collecting information, processing and showing personalized results to the user. Some most common examples of the recommendation system are: YouTube & Facebook news feed or friend suggestion, LinkedIn job suggestion or suggested TV series or movies on platform like HULU or Netflix.

As stated, that a recommendation system works on the data from the user e.g. their search pattern, visited places or pages etc. In many situations, system is unaware of this information. That makes it very difficult for the system to predict the data for a new user. This is called cold approach problem.



1.2 Cold Start Issue in the recommendation system

Any recommendation system works on the data collected from the user e.g. their search pattern, visited places or pages etc. In many situations, system is unaware of this information; e.g. A new user visiting an ecommerce website, new person joining a job portal.

That makes it very difficult for the system to predict the data for a new user. This is called cold approach problem.

1.3 Purpose

The whole project revolves around the cold start problem in Collaborative filtering-based recommendation systems. First, a basic collaborative filtering-based recommendation system was built and then one with the variations to address the cold start problem. The purpose of this project is not to improvise the existing collaborative filtering-based recommendation system but to see if the available information from the user can be used to ease the recommendation process.

Also, I did a literature survey on some recent papers in the reputed journals. This was done to keep up with the recent research on this area and identify some where this project can be extended in the future.

2. Basics of a Recommendation System

As mentioned in the introduction that a recommendation system is used to ease the problem of information overload, information overload is when user finds it really difficult to find the right thing. The information overload problem is addressed by either one of them:

1. Information Retrieval
2. Information Filtering

Therefore, the solution to the information overload problem can be dealt with one of them.

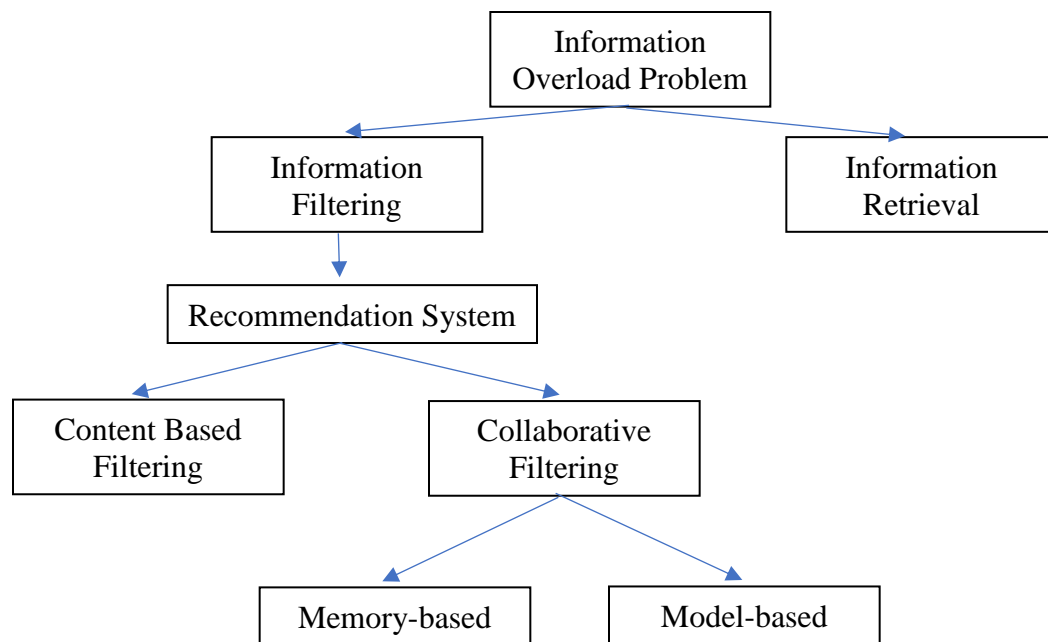


Figure: Information Overload Problem

Information Retrieval Systems directly ask users for the information they need to narrow down the results while the Information Filtering Systems learns it from the user's profile, e.g. his past search results, age etc.

As the collaborative filtering is a part of the Information Filtering, so the focus of this report would be on Information Filtering Side.

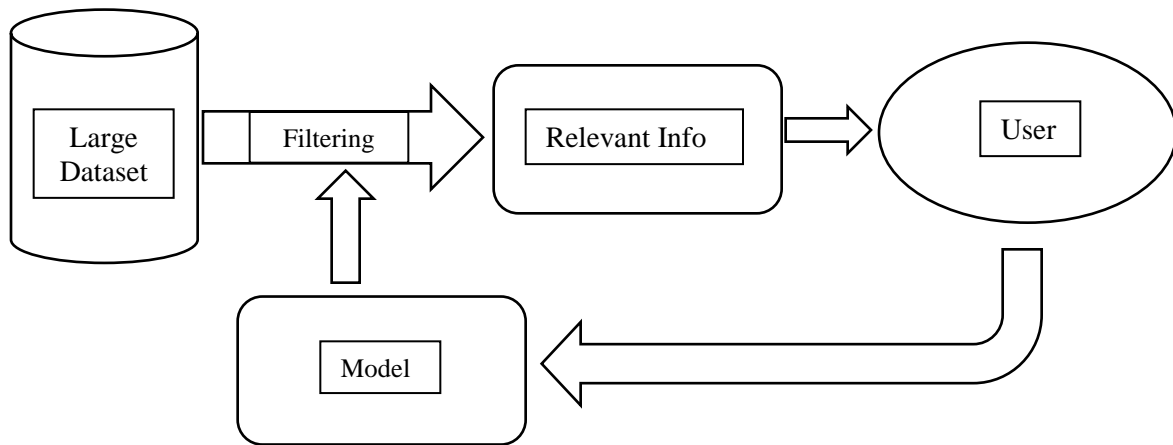


Figure: Information Filtering

The first recommendation system was built by Goldberg in 1992, this recommendation system was based on solving the problem of numerous amounts of emails presented in the inbox. This system was a type of collaborative filtering where users rate their emails that were sent to the user's inbox.

Since the inception of the first recommendation the basic idea has remained the same to offload the not required information, but the overall architecture of the system and applications have witnessed innumerable changes.

As the number of applications that use the recommendation system have been growing, the need for a precise and accurate recommendation has gone up too. Companies and the universities who deployed this technique into their system have reported that these techniques are really very beneficial.

A recommendation system can be summarized with the help of figure 1.1

2.1 Basic Recommendation Systems developed in past years

1. **Collaborative Filtering:** In this type of system, the recommendations are based on the ratings given by the user and similarity between them.

2. **Content Based Filtering:** It is based on the closeness between the products.
3. **Demographic Filtering:** This type of system uses the user's demographic information like age, occupation and location etc.
4. **Social Filtering:** In this system, user's social network's information is used to make the recommendation. Suppose if I recently bought something which is popular among my age group it would be suggested to my friend too.
5. **Hybrid Filtering:** This can simply be a combination of one or more of the systems described above.

As the whole project is based on the Collaborative filtering, therefore the next section has been devoted to explore a little about them.

2.2 Content-Based Filtering

In the type of system, the items that have already been rated by different users are considered and similarity between them is computed. This makes the basis of recommendation.

As mentioned earlier in the report Information retrieval and Information filtering. In this user is not trying to make a search or query, the system has to learn about the user depending on his past behavior and selection and present the best recommendation. This technique is also known as cognitive filtering because the recommendation is based on implicitly learning the user behavior then directly asking for the required information. This type of filtering has taken some concepts from the information retrieval.

In the information filtering, a vector is generated from the user's past behavior, this vector represents the user profile to the system. Depending on this profile vector, the similarity with different products is calculated. Basically, this vector represents a user's behavior in a way that can be interpreted by the computer.

2.2 Collaborative Filtering

This type of system relies on the user's ratings, which is used to calculate the similarity. It is one of the mostly used technique because it is both easier to interpret and implement and gives good results.

Any collaborative Filtering based system focuses on two points: how to calculate similarity vector and how to use it to make the predictions.

In simple words this technique depends on the model that is generated from the user's profile. This model can be based on the profile of a single user or a group/collection of users' profiles as the name also suggests.

Collaborative filtering-based algorithms can be divided into two parts:

1. **Model Based Filtering:** In this type of filtering the system doesn't have to go through the whole dataset instead a predictive dataset is generated from the available data, which is used whenever we have to make a recommendation.

Example: Item Based Filtering Algorithm

2. **Memory Based Filtering:** In this type of filtering algorithm the system has to go through the whole dataset, hence it takes up a lot of resources to build and running it at runtime requires a lot of processing power.

Example: User Based recommendation system

2.2.1 Item based Collaborative Filtering

As mentioned above Item based Collaborative Filtering, also referred as IBCF is an example of Model Based Filtering technique.

In this method, a similarity matrix is required to be calculated. While making the recommendation, relationship between items is taken into consideration. Examples of the similarity measures: Cosine Similarity, Pearson Similarity.

In the layman terms, we maintain a matrix of where we save the k items which are similar to this product. To make the recommendation more accurate, the rating given by the similar user is preferred.

For the prediction, we just calculate the similarity score with one of the similarity measures. And the one with the highest score is considered.

2.2.2 User based Collaborative Filtering

As described above, this comes in the category of the memory-based filtering technique, it is also referred as UBCF. The whole behind this is based on that the users who have similar interests in the items would give similar rating to the products.

Therefore, If we can find the users who have given similar rating to some products would have similar choice. This can be used to make predictions.

2.2.3 k Nearest Neighbor Algorithm

This algorithm has been widely used in the algorithm, that is why it requires a mention in the report. It is both easy to understand and gives great performance. This algorithm divides the data into classes, the similarity between the point is simply the distance between them.

While making the prediction, this algorithm would find out the k most close point. It would find out to which class most of the points belong, the test item is said to belong to that class.

Because the kNN doesn't generalize, therefore this algorithm has to maintain the whole dataset for the testing phase. This is one of the reasons why the training phase is quite fast while testing is little slow.

In simple terms, we pick a constant k. This denotes how many points would be there in a class. The k points that are very similar to each other are classified into a class. kNN can be simply interpreted with the help of a cartesian plot.

2.2.4 Hybrid Filtering Techniques

In the hybrid filtering technique, we combine two or more of the algorithms. This can definitely provide great results but comes with a cost of more processing.

Example: Content based filtering based on demographics.

3. Literature Survey

This part of the report has been dedicated to the literature survey of existing research work which deals with the cold start problem.

Cold start problem has been one of the most challenging problem while building a recommendation system. Imagine opening a youtube page on incognito mode, everything looks so weird. If you manage to click one of those links, your whole feed gets messed up. Hence, it is very important to address the cold start issue.

Cold start problem can be divided into two parts, when a new user logs into the system and when a new product is introduced in the system.

When a new user logs into the system, very little information is present about that user so the predictions can be very inaccurate. Many of the recommendation systems deal with this problem by asking to rate technique, simply they ask the user to rate some of the products. While in case of the new product the collaborative filtering is quite hard, because the collaborative filtering is based on the collaborative filtering and it's not present for a new product. Content based filtering works quite better in this case.

In the coming sections, I have tried to address some of the existing research work published in the reputed journals

Using Demographic Information to Reduce the New User Problem in Recommender Systems

In this paper, authors have tried to use the demographics information present in the MovieLens 100K dataset. I am also using the same dataset in this project. This dataset includes the 100 thousand ratings from 963 users. All these ratings are for 1682 movies. Ratings are on a scale of 5.

Demographic information has the following format:
user id | age | gender | occupation | zipcode

User ratings are present in this format:
user id | item id | rating | timestamp

The dataset is then further divided into test and training datasets, a model with k different groups/clusters is created. This is done with the training dataset while in training process. This model takes the user demographics into consideration while making predictions.

Given plot shows the predicted ratings for a given user based on the information present in that database.

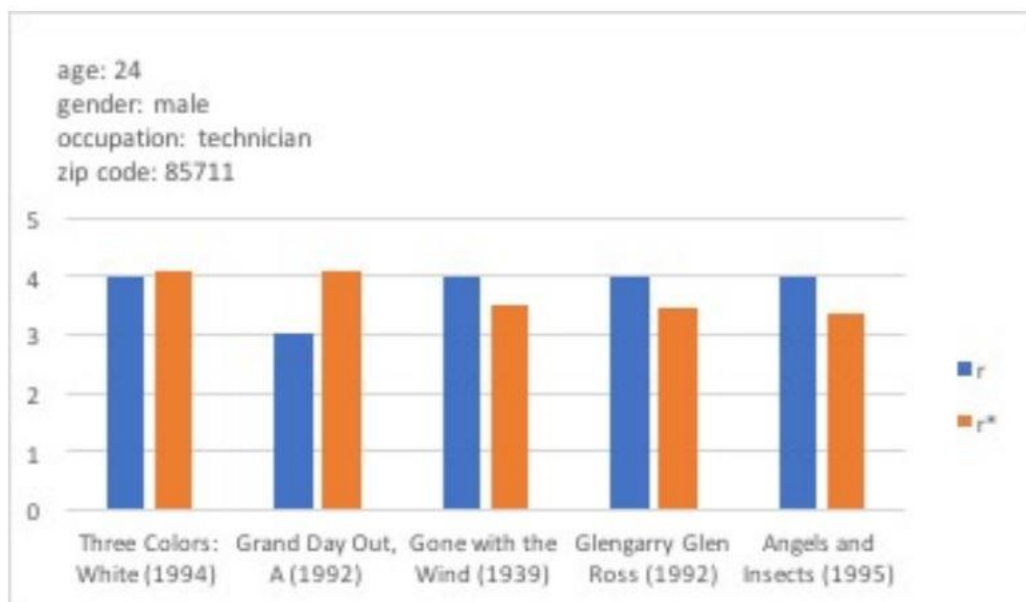


Figure: r^* denotes the predicted rating using the given information in the plot that is occupation, gender, age and zip code.

This plot shows that there exists a relationship between the demographic data and movie rating.

Collaborative Filtering Enhanced By Demographic Correlation

In this paper, the author has combined some of the existing algorithms with the demographic data. This paper introduced two algorithms named U-Demographic and I-Demographic. Where “U” in the first name comes from the user-based collaborative filtering while the “I” in the second name comes from

the user-based collaborative filtering. The author in this paper has further divided the dataset into different age groups. It gives a better representation of the data and results are more accurate.

feature #	feature contents	comments
1	age ≤ 18	<ul style="list-style-type: none"> • each user belongs to a single age group, • the corresponding slot takes value 1 (true) • the rest of the features remain 0 (false)
2	$18 < \text{age} \leq 29$	
3	$29 < \text{age} \leq 49$	
4	age > 49	
5	male	<ul style="list-style-type: none"> • the slot describing the user gender is 1 • the other slot takes a value of 0
6	female	
7-27	occupation	<ul style="list-style-type: none"> • a single slot describing the user occupation is 1 • the rest of the slots remain 0

Figure: Demographic information used by the author.

A new user is molded into a vector form filling up all the fields in the above chart in binary format. The results show the performance is much better than the base model.

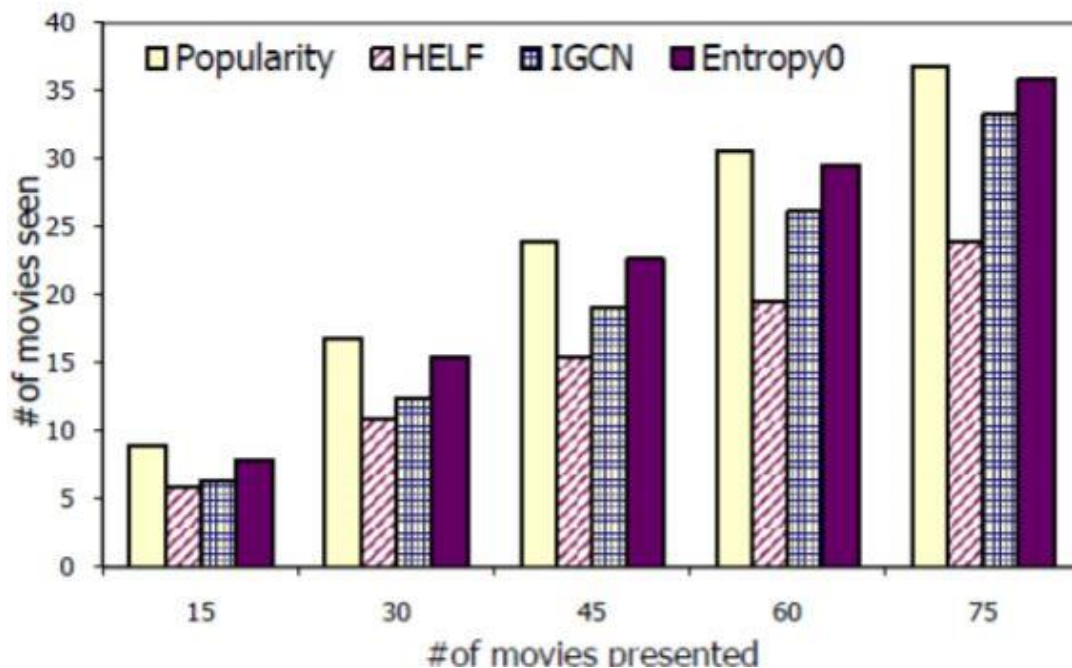
Cold-start Problem in Collaborative Recommender Systems: Efficient Methods Based on Ask-to-rate Technique

This paper tries to enhance the old ask to rate technique, depending on the rating first the similarity with a user in the dataset is calculate. Then from the active user list products are recommended.

Learning Preferences of New Users in Recommender Systems: An Information Theoretic Approach

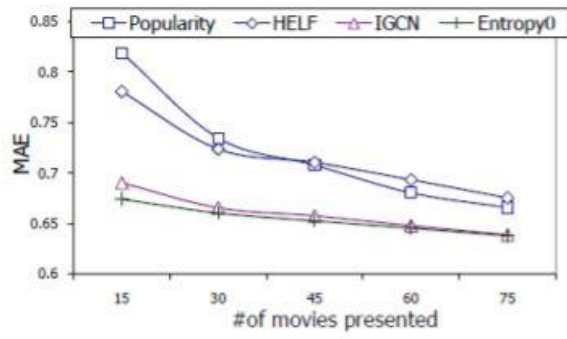
In this paper, author has tried to address cold start problem by calculating the effectiveness of different methods that are based on information theory.

Some techniques were developed to extract the information from the new users to learn about them and make recommendations.

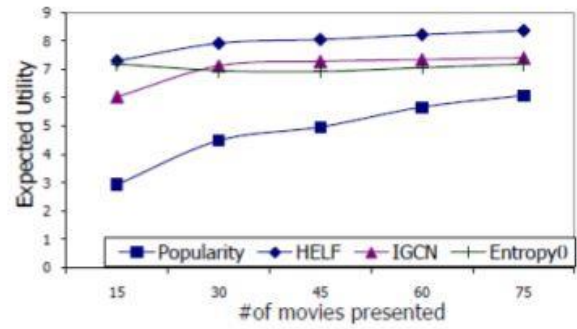


As observed in the previous section, the author also realized problem with the entropy method, so he proposed a new variation named as entropy0 and HELF.

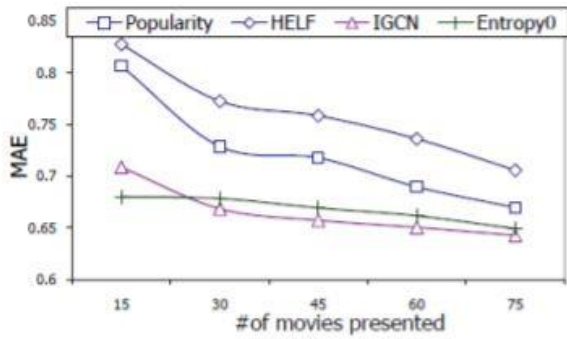
Author worked on 4 methods: HELF, Popularity, IGCN and Entropy0. It comes out be true that the popularity method is the most effective one.



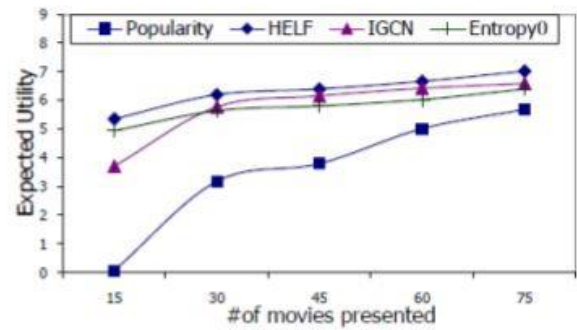
(a)



(b)



(c)



(d)

4. System Architecture

This system has been implemented to exploit most of the information present in the dataset. There can be two form of this system, one interactive and other non-interactive system, where the actual people can be asked to sit in front of the screen and try their hands on with the system but later, in the non-interactive system, simply available information in the dataset is used.

Input: MovieLens 100k dataset

Output: Recommended movies:

- 1) On random basis: This serves as the base model. Some random movies are shown to the user, then collaborative filtering is done. Every other technique is compared to it.
- 2) Movies for rating based on the demographic score: Movies presented based on the demographic information e.g. age, gender, location. Then collaborative filtering is done.
- 3) Movies for rating based on the entropy0 score: Movies presented based on the entropy0 score, then collaborative filtering is done.
- 4) Combination of Demographic and entropy0 score: Movies are displayed in considering this information, then collaborative filtering is done.

kNN algorithm has been used to implement this system, also the bias (average rating given by the user is subtracted). This is done because, some users are too generous while giving ratings while some are too harsh. Subtracting the average rating value before doing the calculating helps in removing the bias. The average is added back at the end after the filtering is done.

4.1 Dataset

As mentioned in the earlier chapter, MovieLens 100k dataset has been used for this project. This data contains a lot of information, so the required information has been taken from it.

The first needed information is user demographics. That is present in u.user file. The actual format looks like this:

User id | age | gender | occupation | Zipcode

It was converted to this format, because the zip code information is simply not required.

User id | age | gender | occupation

For the precise calculation, the same format has been followed as the [4]. All this information is converted to a vector form.

1. Age has been divided into the following group: 0-18, 19-24, 25-30, 31-40, 41-50, 51-60, 61-70, 71-100
2. 0/1 corresponding to M/F for gender identification.
3. Occupation has 21 categories in the dataset.

Age vector is a binary row vector of size $[1*8]$

Gender is another binary row vector of size $[1*2]$

Occupation is a binary row vector of size $[1*21]$

A completed demographic feature looks like: [Age | Gender | Occupation] is of size $[1*31]$

Example:

1. A 10 year old guy would look like this in the given format:
[1, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0]
2. A 23 year old female scientist would look like this:
[0, 1, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0]

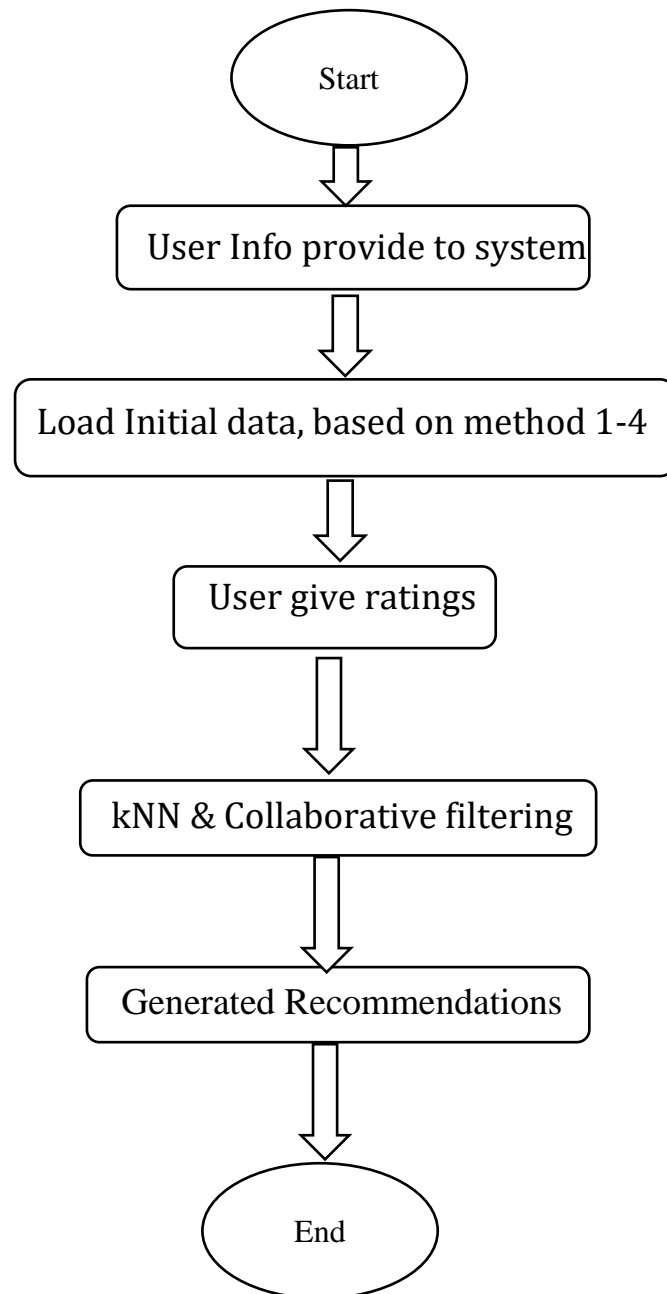
The second required information in the dataset is the user and their corresponding ratings data.

This information is present in u.data file. It look like this.

User id | movie id | rating value

Information related to the movie has been given in the u.item file but only movie name is of use when we use interactive system, otherwise in case of non-interactive system this dataset is of no use.

4.2 Architecture



5. Results and Conclusion

After the system was built, it was the time to test it on the real people. Unfortunately, this was the most challenging part of this research. The dataset includes a lot of diverse dataset but, I could only test a limited part of it as my friend circle is very limited by in the age and profession.

I asked a total of 20 people, 13 of whom were male and 7 were female. 4 different scripts were run and they were asked to choose movies on each of them based on their preference. No of the users were told about which script they are running

Type 1 corresponds to the system base on just the demographic, type 2 on entropy0 score, type3 is a hybrid model of 1&2, and type 4 is the one that shows random movies to the users.

The users were asked to rate their experience with all the 4 scripts, a score of 4 represents the most preferred one, 1 represents the least preferred. The raw data can be seen in the table given below.

User Demographics			User Preference Score			
Age	Gender	Occupation	Type 1	Type 2	Type 3	Type 4
30	M	Engineer	3	4	1	2
22	F	Student	1	3	4	2
38	M	Technician	4	3	2	1
18	M	Student	4	3	1	2
20	M	Artist	4	3	2	1
24	M	Engineer (Just Graduated)	3	4	1	2
43	F	homemaker	3	4	2	1
39	F	homemaker	1	4	2	3
35	M	Other (Gardner)	4	1	2	3
23	F	Student	3	4	1	2
65	M	retired	4	1	2	3

27	M	Programmer	3	4	2	1
31	M	Student (Economics)	2	4	3	1
26	M	Scientist (Data)	3	4	1	2
29	M	Engineer (Graduated)	1	4	2	3
26	F	Educator (Elem. School)	3	2	1	4
23	F	Programmer	3	4	2	1
28	M	Student (Economics)	4	2	3	1
40	M	Other (Mechanic)	3	4	2	1
24	F	Educator	4	2	3	1
Sum			60	64	39	37
Average			3	3.2	1.95	1.85

The entropy0 based system came out to be the most preferred one, followed by demographic based, followed by the hybrid model. The worst experience was observed with the system that throws up the random movies to the user.

Future work

As it is correctly visible that there is a correlation between the demographic information and the choices people make. Therefore, this project can be further expanded for other areas like job portals, e-commerce websites for different kind of products. Unfortunately, it couldn't be done because of the absence of the data, this dataset was just limited to the movies.

Also, if we could introduce more demographic information in the dataset e.g. religion, language, ethnicity information. The performance could be further improved.

References

1. <https://machinelearningmastery.com/practical-machine-learningproblems/>
2. <https://www.ritchieng.com/machine-learning-project-customersegments/>
3. Callvik, Johan, Liu, Alva, (2017) “Using Demographic Information to Reduce the New User Problem in Recommender Systems.” Kth Royal Institute of Technology School Of Computer Science And Communication
4. Manolis Vozalis, Konstantinos G Margaritis, (2004). “Collaborative Filtering Enhanced By Demographic Correlation.” AIAI symposium on professional practice in AI, of the 18th world computer congress.
5. MH Nadimi-Shahraki, M Bahadorpour, (2014). “Cold-start Problem in Collaborative Recommender Systems: Efficient Methods Based on Ask-torate Technique”. CIT. Journal of Computing and Information Technology 22 (2), 105-113
6. Rashid, A.M., Karypis, G., Riedl, J. (2002) “Learning preferences of new users in recommender systems: an information theoretic approach” 127– 134. ACM Press