# An Approach for Friend Recommendation Based on Selected Attributes

Rachana Raghuwanshi, and Gend Lal Prajapati, *Member, IAENG*

*Abstract*—**Recommendation algorithms are becoming vital for their use in social networking applications. Users' lifestyle plays the important role in identifying users' nature and in explicitly rating to represent their interest. Lifestyle based friend recommendation can be an innovative approach to find online friends and to make a group as per their interest. This paper aims to recommend friends of a user based on their selected attributes. In this framework, we have collected real data of 200 users by utilizing the Google form on the basis of 52 different attributes. We use lifestyle extraction and quantification approaches for the characteristic extraction. Afterward, K-means clustering approach is applied to make a group of similar users based on closeness in their lifestyle. We also apply collaborative filtering based similarity indexing approach to the relevant cluster to find more close and accurate results. The complete work is implemented in JAVA technology and tested on the basis of precision, recall and F-score parameters.**

*Index Terms*—**feature extraction, friend recommendation, K-means clustering, quantification, similarity index**

## I. INTRODUCTION

PREDOMINATING friend recommendation framework prompts friends to users contingent upon the criteria, which is friends of friends or common friends. Which may not be most right or valuable for picking friends in real life. In this paper, a linguistic based friend recommendation framework for interpersonal organization is utilized, which counsel friends to users based on their life styles rather than common friends or social graph criteria. The framework proposes a closeness factor to quantify the resemblance of ways of life amongst users, and computes users' impact with the help of clustering algorithm. After getting a query, the framework proposes a list of users as indicated by corresponding component to the query. Here, the rules to cluster people together use different parameters including habits/lifestyles, tastes, blood group, location, qualification, research interest, and profession.

Friendbook [1] is a semantic based friend recommendation framework, taking advantages of sensor-rich Smartphones. Friendbook finds existing styles of users from user-centric sensor records. It measures the closeness of life patterns among users, and recommends friends to users if their lifestyles have high similarity. It has suggested the similarity metric to measure the similarity of life styles between users impact in terms of lifestyles by using a graph. Friendbook returns a list of people with highest rating to the query users. It integrates a mechanism to further enhance the recommendation accuracy. It proposes that recommendation systems can be classified into two Regions of centers: link recommendation and object recommendation. First one is person to person communication (e.g., Facebook [9], LinkedIn), and second one is item or objects suggestion (e.g., Amazon, Netflix). The proposed work in the paper [1] is focused on the first one (link recommendation). More study about Friendbook can be seen in [7] [8].

System by Bian and Holttzman have presented Collaborative filtering friend recommendation system depending on personality matching [3]. Kwon and Kim [6] provided a friend recommendation system which is using physical and social context.

Video suggestion framework being used at YouTube, the world most prominent online video group. The framework prescribes customized sets of video to user based on their activity on the site. Video recommendation system, which delivers personalized sets of videos to sign in users based on their previous activities on the YouTube site can be seen in [10]. Amazon utilizes recommendation based on item to item collaborative filtering procedure. The algorithm produces suggestions based on a few customers who are most likely users. It can measure the similarity of two customers and measures the cosine of the angle between the vectors [4]. Netflix recommendation system gives video on request in UK. Netflix has this goal of trying to make suggestions of which movies or TV shows might be of interest to the person who has come to the site [5].

## II. PROBLEM STATEMENT

Friend recommendation is vital and requiring use cases of social networking portals. It becomes more significant when an organization considers it for marketing, and identifying peer users matching. Lifestyles reflect the nature of users along with their interests, and can be used to predict users based on their nature [2]. Study of existing solutions in literatures concludes that user classification, clustering, and similarity matching algorithms are the most common ways of user mapping and recommendation purpose. All such techniques consider user transaction and action to diagnose the relationship and similarity analysis. Some of them considered lifestyle as the lead reason for user matching and

Manuscript received March 07, 2017; revised April 11, 2017.
Rachana Raghuwanshi, is currently enrolled with the Department of Computer Engineering, Institute of Engineering & Technology, Devi Ahilya University, Khandwa Road, Indore - 452001 INDIA (e-mail: rachana.raghuwanshi1990@gmail.com).
Gend Lal Prajapati, is with the Department of Computer Engineering, Institute of Engineering & Technology, Devi Ahilya University, Khandwa Road, Indore - 452001 INDIA (phone: 0731-2366800; fax: 0731- 2764385; e-mail: glprajapati1@gmail.com).

recommendation purpose. A Lifestyles based friend recommendation system is still need to advance. And hence in this paper, we use additional attributes of lifestyles to test the efficiency of a friend recommendation system. We have collected the real data from various level of users. We also use K-means clustering algorithm for exact matching. K-means is very quick, versatile, robust and easier to implement and deploy. Here K-means clustering algorithm is used with the recommendation feature.

### III. PROPOSED METHODOLOGY

Step 1: Data Collection and Preparation:
   User Information is stored in the database.
Step 2: Parsing:
  a) Data cleaning: We remove null characters (like blank space) and stop word.
  b) Lemmatization: We use Stanford lemmatizer library.
Step 3: Determining Similarity:
  a) Quantification: We use Quantification module to exact the lifestyle of a user on the basis of matching formula.
  b) K-means clustering: We apply K-means clustering algorithm that show the list of friends having similar interest.
Step 4: Recommendation:
  a) Similarity weight calculation: We calculate the highest similarity factor on the basis of quantification approach, and on the basis of similarity index for recommending the friends with highest similarity index.
  b) Threshold Filter: Threshold based clustering automatically filters the recommendation according to the value of the threshold. Which automatically filters the recommendations according to the value of the threshold.
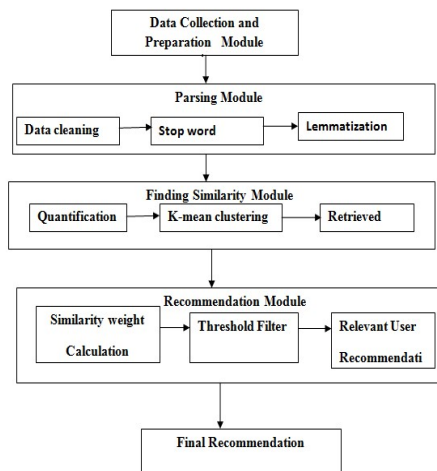   The methodology is also sketched in Fig. 1.



Fig. 1. System Architecture of Recommendation System

### IV. DATA COLLECTION AND PREPARATION MODULE

Data collection can be done into two ways: either from a direct user or from the existing system. Interview, questionnaires, the survey can be a great route for essential information gathering. Subsequently, the auxiliary information can be acquired from a prior framework or dataset. We considers primary data collection as the prime source for evaluating purpose. In data collection module, we have gathered information by utilizing the Google form on the basis 52 attributes more than 200 users. The Collection module gathers life documents from the browser and stores it into an excel file in either semi-organized or organized configuration (database link is given below). Sample of All 52 Attributes are Shown below in Table I. Click on the database link to open Google form and download the excel file of the database. Which contains 52 attributes and various responses. We have collected more than 132 responses by using Google form and more than 70 responses by offline. Data collection is done by college students, faculties, staff members, relatives and friends. Below are some notable points with reference to our database:

1. If a user does not fill the value of some attribute then is consider as a null entry (use data cleaning.)
2. sometimes a user enters the attributes information in capital letters, so we use Stanford lemmatizer to convert it to base format and then cheek
3. When a user fills information like "CSE" in place of "computer science engineering" we use abbreviation rules for matching this information.

A sample entries are shown in Table I of two instances (User Ids) based on 52 attributes. Complete Information is given below on using a database link.

https://docs.google.com/forms/d/1SLvXzMY8J87vsLYEnI-tbxvieK4v7U2PQsVhcPBAgVc/edit#responses

### V. PARSING MODULE

Data cleansing is the way toward distinguishing and revising degenerate or wrong records from a database. Then subsequently, stop words removal and lemmatization has been developed to improve the quality of data source. Stop words are deemed unessential for searching purposes because they occur frequently in the language for which the indexing engine has been tuned. In order to save both space and time, these words are dropped at indexing time and then ignored at search time. Afterward, Tokenization scheme and abbreviation removal have been applied to break complete work into individual values and extract more accurate lifestyle answers. The complete work explores more accurate lifestyle values. The complete data collection form view has been shown in Fig. 2.

### VI. SIMILARITY MODULE

*A. Quantification*

The Lifestyle of users extracted by using the quantification module. We quantify lifestyle by using the closeness factor. To compute the closeness factor we use the following formula:

Closeness factor $= 2*R/(S_1+S_2)$

Where, $R$ is total number of matching words in both the strings, $S_1$ is total words count in string $S_1$, and $S_2$ is total words count in string $S_2$.

TABLE I
User ID with Attributes

| Attributes No. | Attributes Name | User ID 1 | User ID 2 |
|---|---|---|---|
| 1. | Name | Pramila Patidar | Preeti Goyal |
| 2. | Nick Name | Pammi | Preetu |
| 3. | What is your ultimate Aim or Goal | Became a successful teacher | Became successful person |
| 4. | Email | ppatidar@sdbce.ac.in | goyalpreet.02@gmail.com |
| 5. | Address | kasrawaad | Indore |
| 6. | Mobile | 7509765893 | 7415724471 |
| 7. | Date of Birth | 09/11/1987 | 02/08/1989 |
| 8. | Gender | Female | Female |
| 9. | Interested In | Man | Man |
| 10. | Blood Group | O+ | B+ |
| 11. | Hobbies | dancing | leraning new things |
| 12. | Religion | Hindu | Hindu |
| 13 | Add your political view | Gender equality | Gender baised should remove |
| 14. | Add your Religious view | I believe that everything | All are equal |
| 15. | School name and date attended | j.n.v. sanawad and 2003 | Girls Higher sec. school , 2005 |
| 16. | School name and date attended | j.n.v. sanawad and 2005 | Girls Higher sec. school 2007 |
| 17. | Degreessplization and date attended | BE,CSE , 2009 | BE,CSE , 2011 |
| 18. | Specializationand date attended | Computer Science and 2016 | ME , IET DAVV , 2017 |
| 19. | Research Interest | Network Security | Cloud Computing |
| 20. | Work Place | S.D.B.C.E. indore,Assit. Prof. | S.D.B.C.E. indore,Assit. Prof. |
| 21. | Communication Language | English | English |
| 22. | Professional Skill | Information Technology Skill | Java |
| 23. | Extra Activities | high jump,long jump | Drawing |
| 24. | Relationship status | Married | Single |
| 25. | Current city | Indore | Indore |
| 26. | Home Town | kasrawad | Khandwa |
| 27. | Yoga | Yes | No |
| 28. | Regular Exercise | No | Yes |
| 29. | Drinking Habits | No | No |
| 30. | Smoking Habits | No | No |
| 31. | Add a Website | google and youtube | Instagram |
| 32. | Add a Social links | www.google.com | - |
| 33. | Favorite Sports | kabaddi | Football |
| 34. | Favorite Sports Person | sachin tendulkar | Virat Kohali |
| 35. | Favorite Movie | mother india | Hum aapke hai kon |
| 36. | Favorite Song | jag ghumiya tere jaisa na koi | tum pas aaye |
| 37. | Favorite Artiest | Amir khan and salman khan | AR Rahman |
| 38. | Favorite Book | mansarovar | Complete reference of java |
| 39. | Favorite TV Shows | Big Boss | Tarak mehta ka ulta chashma |
| 40. | Favorite Colors | Red | Blue |
| 41. | Favorite Place | Goa | Maharashtra |
| 42. | Favorite Brands | Blue Heven | Lakme |
| 43. | Add a Family and Relationship Member | Mr. Mahendra Patidar | Nidhi Goyal , Sister |
| 44. | Favorite Quotes | Give respect and take respect | You can never get enough. |
| 45. | Favorite Vacation Destination | pachamadi | Goa |
| 46. | Favorite Foods | Burgar | Italian |
| 47. | Online Shopping | Yes | No |
| 48. | Type of clothing | Indian Traditional | Indian Traditional |
| 49. | Favorite Clothing Brand | Viva | Viva |
| 50. | Favorite Watch | Fastrack | Fastrack |
| 51. | Favorite travelling way | Aeroplane | By Road |
| 52. | Are you Supporting fordemonetisation | Yes | Yes |

*B. K-means clustering*

Clustering is an approach to classifying all elements in such a way that every similar element should reside into the single group based on their similarity. Subsequently, it also resides irrelevant elements into another group based on their similarity value and maximum cluster size. Here, K-means clustering approach has been used to construct a group of similar users based on lifestyle similarity. It is one of the simplest unsupervised learning algorithms which simplifies the work of mining by classifying the similar elements in a cluster using the k-cancroids parameter. It calculates a distance between each element to evaluate similarity and reside them into a single cluster by comparing with the k-centroid parameter. The Euclidean distance function measures the distance from point A to point B. The equation for this distance between a point P ($y_1$, $y_2$, ..., $y_n$) and point Q ($z_1$, $z_2$, ..., $z_n$) are:

$$D = \sqrt{\sum_{j=1}^{n} \left( y_j - z_j \right)^2}$$

A snipping of this work is shown in Fig. 3.

| Id | Name | Nick Name | What is your ultimate Aim or Goal | Email | Address | Type of cloth | Favorite C | Favorite V | Favorite t | Are you Su |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | Pramila patidar | pammi | Became a successful teacher | ppatidar@sdbce.ac.in | kasrawaa | Indian Traditi | viva and d | fastrack | Aeroplane | Yes |
| 2 | Preeti Goyal | preetu | Became successful person in the life | goyalpreet.02@gmail.com | Indore | Indian Traditi | viva | fastrack | By Road | Yes |
| 3 | poornima tiwari | punnu | To be a successful person in life | poonimatiwari1@gmail.com | garden no | Indian Traditi | frock | maxima | By Road | Yes |
| 4 | Rishabh Vishwakarma | Rishu | Software engineer | rishabhvishwakarma758@gmai | Goramach | Morden | viva | sonata | By Railwa | Yes |
| 5 | nisha salawat | nili | to become a successful web designer | nishasalawat1@gmail.com | 9/36 imli k | Indian Traditional | | sonata | By Ship | Yes |
| 6 | Anand Singh Thakur | Addy | To be a successful person in life. | anand10bundela@gmail.com | Near mair | Indian Traditi | adidas | LED watch | By Road | Yes |
| 7 | Nandita Senapat | Jiya | placement in good company | kishi.senapat97@gmail.com | 87/2E Raje | Morden | viva and d | fastrack | By Road | Yes |
| 8 | Priyanka patidar | piyu | placement in a good company | patidarpiyu2686@gmail.com | 29/6 princ | Morden | viva and d | fastrack | By Road | Yes |
| 9 | Rahul Gupta | Rahul | CTO | rahul.gupta@serosoft.in | 136 Phadn | Western | DENIM | fastrack | By Road | Yes |
| 10 | Rashmi raghuwanshi | Rashmi | Bank manager | rashmi.raghuwanshi1992@gma | Ambedka | Indian Traditi | Max | Fastrack | Aeroplane | Yes |
| 11 | devendra badoniya | rahul | business man | devbadoniya@gmail.com | 142 sita na | Western | raymonds | sonata | By Road | Yes |
| 12 | vishwanath vishwakarm | Vishwa | work as a Network Engineer in Repute | vishwa.vishwakarma84@gmail. | 281/4 Neh | Western | Raymonds | fastrack | By Road | Yes |
| 13 | Rahul patidar | | Environment engineer | patidarrahul03@gmail.com | 1816- D su | Western | Nike | Fastrack | Aeroplane | NO |
| 14 | Sumit Tiwari | Sumit | Entrepreneur in MANUFACTURING AN | tiwari.sumit59@gmail.com | 7/9 kamal | Western | us Polo | fastrack | By Road | Yes |
| 15 | Lalit yadav | | To earn money without harming any or | yadavlalit20@gmail.com | 148 nema | Morden | DENIM | CASIO | By Road | Yes |
| 16 | Rohit soneja | Aalu | Nthing | rohitaalu1419@gmail.com | Savarkar | Morden | DENIM | Rs | By Road | Yes |
| 17 | Damini gupta | | Bank manager | daminigupta1034@gmail.com | Behlot by | Morden | DENIM | fastrack | By Road | Yes |
| 18 | Radhika | Ranu | Any government job | raghuwanshiradhika1489@gma | Kalabag ga | Indian Traditi | No | Fastrack | By Railwa | Yes |
| 19 | Rajesh singh yadav | Raju | To transform my dreams in to reality | rajeshsinghyadav007@gmail.co | 3 dr abani | Western | DENIM | fastrack | By Railwa | NO |
| 20 | Raksha Raghuwanshi | Mishti | Bank manager | raksharakshraghuwanshi1995@ | Behlot by | Indian Traditi | Tashvi | Sonata | By Railwa | Yes |
| 21 | vaidehi dave | vaidu / sho | Software engineer | vaidehidave66@gmail.com | 14 kunwar | Morden | levis | sonata | By Railwa | NO |
| 22 | soumil tiwari | soum | respective position in society | soumil489@gmail.com | delhi | Indian Traditi | DENIM | CASIO | By Road | Yes |
| 197 | vinay choudhary | | being giving | choudharyvinay@gmail.com | bina | Morden | guess | fastrack | By Road | Yes |
| 198 | sachin sharma | sachu | achieving self knowledge | rstsachine@gmail.com | sironj | Western | Peter Eng | Rolex | By Ship | Yes |
| 200 | preeti yadev | daisy | avoiding harm | preeto_ya@gmail.com | umaria | Western | Peter Eng | fastrack | By Road | Yes |

Fig. 2. Data Collection by Using Google Form



Fig. 3. Suggested friends based on K-means

## VII. RECOMMENDATION MODULE

### A. Similarity weight calculation

The most significant challenge in this section was to map the lifestyle of each user and convert them into quantify figures. Here, quantification technique has been implemented to convert all similar values into matching score. Afterwards, output has been forwarded to clustering module for cluster making. This complete phenomenon generates the most relevant users based on similarity distance based on following formula.

Similarity score = total of all similarity matching weight

UrefID = Reference User ID for recommendation

$Sw_{[Total]} = Sw_1 + Sw_2 + Sw_2 + ...Sw_n$

Here, $Sw_{[Total]}$ = total similarity weight of all attributes with respect to UrefID

$Sw_1$ = similarity weight for attribute 1

$Sw_2$ = similarity weight for attribute 2

$Sw_3$ = similarity weight for attribute 3

$Sw$ = similarity weight

n = total number of attributes

The recommendation system is a method of filtering that use rating, similarity score or preference score to predict the frequency of item and elements. Recommendation systems have become increasingly popular in recent years due to a wide area of applications and use into movies, books, articles etc. prediction and suggestions. Here, a customized recommendation algorithm has been developed and append with K-means clustering algorithm to provide more accurate and relevant solution. Here, recommended cluster is used as input data source and similarity score has been calculated. Similarity score represents the total lifestyle closeness of each user with desired user lifestyle. In simple words, high similarity score represents more lifestyle closeness in comparison with users having low similarity score. Similarity threshold value has been used to filter out the retrieved User ID and recommends most close user references. The high threshold would represent filter with high strength and it will recommend users with most close values. Subsequently, low threshold represents low filter strength and high numbers of users.

## VIII. Experiment Analysis

A java based recommendation tool has been developed to implement the proposed solution. Proposed implementation view has been classified into four modules. Initially, all collected data has been exported into .CSV file format from Google docs and loaded to perform parsing process. Incomplete data removal technique has been implemented for data cleaning purpose. Afterward, Stanford libraries are used for stop word removal, lemmatization and tokenization purpose. This module helps to provide more exact and accurate data source for clustering process. The major challenge during clustering process was a mapping of all user lifestyle into a numeric representation. A self-developed quantification process has been used to convert all content sentiments into numeric figures. Subsequently, K-means clustering approach has been performed to extract similar users based on desired user information. Clustering can help us to retrieve relevant users but can't be considered as recommendation technique. A self-proposed recommendation technique based collaborative filtering has been implemented for suggestion purpose. Here, similarity weight has been considered to evaluate the ranking of a user in similarity index cluster. Similarity weight is the total sum of all weights estimated during quantification process. At last threshold value has been used to filter out all retrieved document and generate most relevant users as a final recommendation. Recall-precision and F-score parameters have been used to measure the performance of proposed solution. These User ID are referred from data collection module. The complete view of implementation is shown in Fig. 4.

## IX. Result Analysis

The complete performance has been evaluated on basis of Recall [Accuracy], Precision and F-Score [Final Score]. It has been observed that clusters acquired from the threshold

based techniques are more separated and compact which indicates good clustering. Initially, ten users have been selected on a random basis for friend recommendation purpose. Selected User IDs like 1, 6, 8, 32, 64, 78, 81, 100, 106 and 115. Afterward, five different threshold values have been selected to obtain results for different filtering strength like 6 (Very Low), 9 (Low), 15 (Medium), 22 (High) and 30 (Very High). Selected User IDs are for Cluster size 2.

*Example:* Let the similarity weight of an User ID is 30, and selected threshold value is 6. Then system will recommend all those users whose similarity weight is between 6 and 30 (includly 6 and 30). Hence more number of users will come in this range. On the otherhand if threshold value is large then less number of users will come in the range, and hence system will recommend less number of users. If we increase threshold value then we are getting less number of recommendations, but shows more strength among users. Also, if we decrease threshold value then we are getting more number of recommendations, but shows less strength among users.
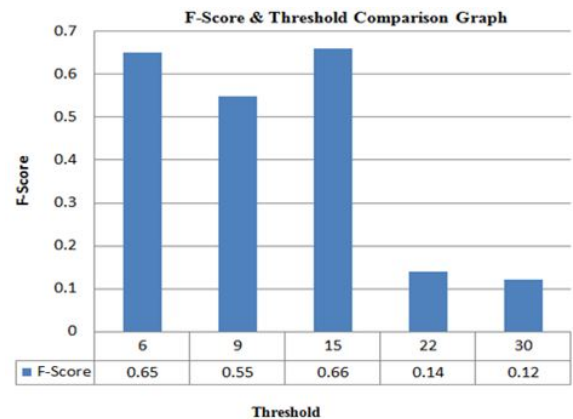


Fig. 4. F-Score Comparison on different threshold

## X. Conclusions

The complete work concludes that a lifestyle based friend recommendation system can add a boost up feature in social networking sites. In this work, a modified version of clustering and filtering approach has been proposed as the hybrid solution for recommendation purpose. Subsequently, a custom quantification and filtering approach have been performed to simplify the recommendation process with accurate performance. The recommendation scheme has been analyzed which are listed in Fig. 4, Fig. 5 and Fig. 6. The followings are the notable remarks:

1. Constant recall with value one has been recorded for 6 and 9 thresholds.
2. Increasing into threshold value decrees the system performance and its lacks gradually with respect to enhancement of threshold value.
3. The variable precision score has been recorded with different users. Maximum 0.86 precision and minimum 0.09 score has been recorded
4. An integrated F-Score has been calculated for each threshold value where minimum 0.12 and maximum 0.66 scores have been recorded.
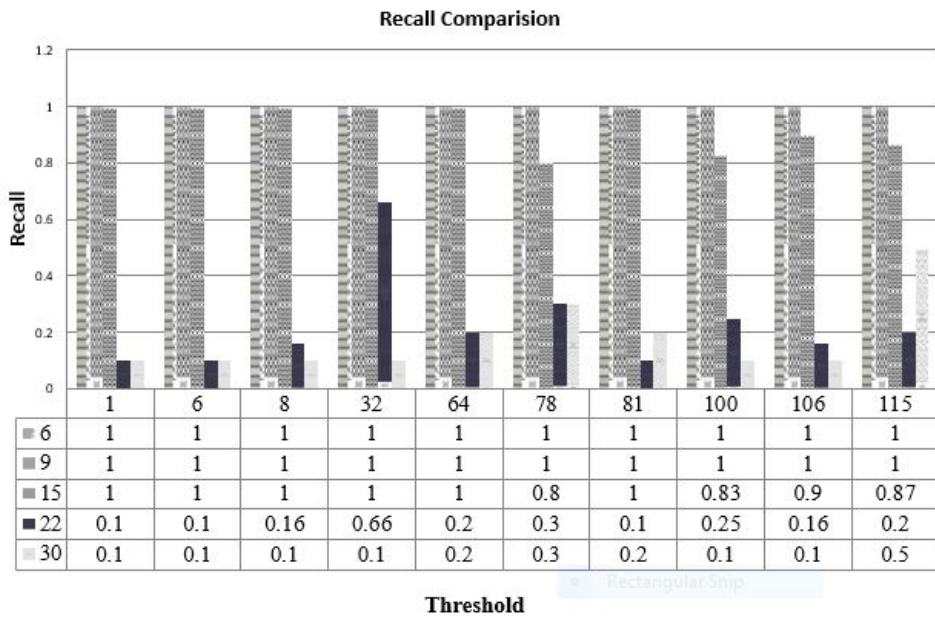
**Recall Comparision**

| Threshold | 1 | 6 | 8 | 32 | 64 | 78 | 81 | 100 | 106 | 115 |
|---|---|---|---|---|---|---|---|---|---|---|
| 6 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| 9 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| 15 | 1 | 1 | 1 | 1 | 1 | 0.8 | 1 | 0.83 | 0.9 | 0.87 |
| 22 | 0.1 | 0.1 | 0.16 | 0.66 | 0.2 | 0.3 | 0.1 | 0.25 | 0.16 | 0.2 |
| 30 | 0.1 | 0.1 | 0.1 | 0.1 | 0.2 | 0.3 | 0.2 | 0.1 | 0.1 | 0.5 |

Fig. 5. Recall Comparison for all users on different threshold

**Precision**

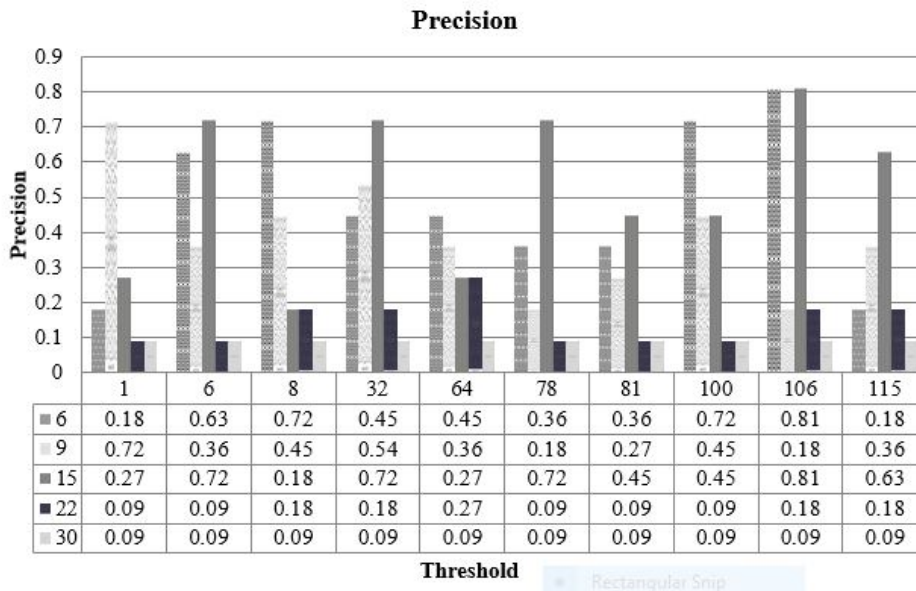| Threshold | 1 | 6 | 8 | 32 | 64 | 78 | 81 | 100 | 106 | 115 |
|---|---|---|---|---|---|---|---|---|---|---|
| 6 | 0.18 | 0.63 | 0.72 | 0.45 | 0.45 | 0.36 | 0.36 | 0.72 | 0.81 | 0.18 |
| 9 | 0.72 | 0.36 | 0.45 | 0.54 | 0.36 | 0.18 | 0.27 | 0.45 | 0.18 | 0.36 |
| 15 | 0.27 | 0.72 | 0.18 | 0.72 | 0.27 | 0.72 | 0.45 | 0.45 | 0.81 | 0.63 |
| 22 | 0.09 | 0.09 | 0.18 | 0.18 | 0.27 | 0.09 | 0.09 | 0.09 | 0.18 | 0.18 |
| 30 | 0.09 | 0.09 | 0.09 | 0.09 | 0.09 | 0.09 | 0.09 | 0.09 | 0.09 | 0.09 |

Fig. 6. Precision Comparison for all users on different threshold

5. It has been observed that clusters acquired from the threshold based technique are more isolated and reduced which demonstrates good clustering.

6. The complete work concludes that proposed solution can be used as the recommendation technique for friend recommendation purpose.

REFERENCES

[1] Zhibo Wang, Hairong Qi,"Friendbook: A Semantic-Based Friend Recommendation System for Social Networks," IEEE Transactions on Mobile Computing, Vol. 14, No. 3, MARCH 2015, pp. 01–14.

[2] Tofik R. Kacchi, Prof. Anil V. Deoranker," friend Recommendation system based on Lifestyles of Users" International Conference on Advances in Electrical, Electronics, Enformation, Communication and Bioinformatics (AEEICB16)2016, pp.01–04.

[3] Bian and H. Holtzman, "Online friend recommendation through Personality matching and collaborative filtering," in Proc. 5th Int. Conf. Mobile Ubiquitous Compute, Syst., Services Technol., 2011, pp.230- 235.

[4] Amazon. (2014) [Online]. Available: http://www.amazon.com

[5] Netflix. (2014) [Online]. Available: https://signup.netflix.com

[6] J. Kwon and S. Kim, "Friend recommendation method using physical and social context," Int. J. Comput. Sci. Netw. Security, vol. 10, no. 11, pp. 116–120, 2010.

[7] Pankaj L. Pingate, S. M. Rokade, "A Survey of Friendbook Recommendation Services", International Journal of Science and Research (IJSR), ISSN (Online): 2319-7064, Vol 3, Issue 11, November 2014.

[8] Z. Wang, C. E. Taylor, Q. Cao, H. Qi, and Z. Wang. Demo: Friendbook: Privacy Preserving Friend Matching based on Shared Interests. In Proc. of ACM SenSys, pp.397–398, 2011.

[9] Facebookstatistics.http://www.digitalbuzzblog.com/Facebook-statistics-stats-facts-2011/

[10] YouTubeAPI Documentation http://code.google.com/apis/youtube/overview.html