

# Clustering and Recommending Services based on ClubCF approach for Big Data Application

G.Uma Mahesh

*Department of Computer Science and Engineering  
Vemu Institute of Technology, P.Kothakota, Andhra Pradesh, India*

G.Lokesh

*Department of Computer Science and Engineering  
Vemu Institute of Technology, P.Kothakota, Andhra Pradesh, India*

**Abstract:** The number of services emerging on the Internet are generating huge amount of data leading to big data. Storing such data using traditional storage approaches is impractical which can be solved using Big Table capable of storing number of services in the form of multi dimensional sorted map again searching for a services and to recommend it to the new users requires large computations. In the present work these problems are solved by using the Clustering based Collaborative Filtering (ClubCF) approach and Mash Up data set with 6888 services along with their description and their functionality is considered for clustering with the help of agglomerative hierarchical clustering algorithm.

**Key Words:** Clustering, Collaborative Filtering, Recommendations, Big Data, Mash Up data set.

## I. INTRODUCTION

The Collaborative filtering technique is an important approach to develop recommender systems that helps in predicting the interests of active users based on the preferences or interests of previous users provided in the form of ratings and reviews.

Traditional collaborative filtering techniques include Item based Collaborative Filtering and User based Collaborative Filtering techniques [1], both the CF techniques suffer from the data sparsity problem because the users rate only few set of services if further the services increase it will be a problem due to the limited users. In Item based Collaborative Filtering [2] the similarity between the services will be calculated and recommended to the new user with the most rated services. In user based Collaborative Filtering technique the similarity between different users is calculated and recommend to the new user with the service that most preferred by the previous users. The Collaborative Filtering techniques will have some challenges for the Big Data applications such as: performance decreases when the data was sparse; cannot provide recommendations to new users and new items because of limited users who will rate the limited number of items; limited scalability for large data sets.

The problems occurred during traditional Collaborative Filtering techniques can be solved by using ClubCF approach which helps in decreasing the number of services that needs to be processed and reduces the online computation time. The ClubCF approach consists of two steps; Clustering and the Collaborative Filtering, during the clustering step different services can be clustered or grouped in to different clusters by calculating their similarities between them and during the Collaborative Filtering step helps to recommend the services to the users by calculating the rating similarities between different users. The benefits of using Club CF approach are; better addresses the sparsity, scalability and other problems; to handle large amount of services in effective manner; to decrease the number of services needed to be processed; improve prediction performance and improve efficiency of recommendations.

## II. SERVICE BIG TABLE

Big table was designed as a distributed storage system aims at storing structured data having ability to vary Peta bytes of information [3]. Some examples of services that are stored by Google people in big table are like Google Maps, Google Book History, Google Earth, Google Code History, You Tube, Gmail etc. The benefits of Big Table includes: distributed, Sparsity, Multidimensional sorted map, Scalability, easy to add new services in an effective manner.

The Google Big Table was indexed by the Row Key, Column Key and a time stamp where each service id or service number can be considered as row key and the column keys which deals with the service grouped together represents as column families, time stamp will focus on The time when the service actually recorded in the big table . The Table 1 will describe an example of how services are stored in big table which represents the 4 wheels route map service indexed by service id or row key as s1 and column families includes description, functionality and ratings provided by appropriate users after using the service along with the time stamps.

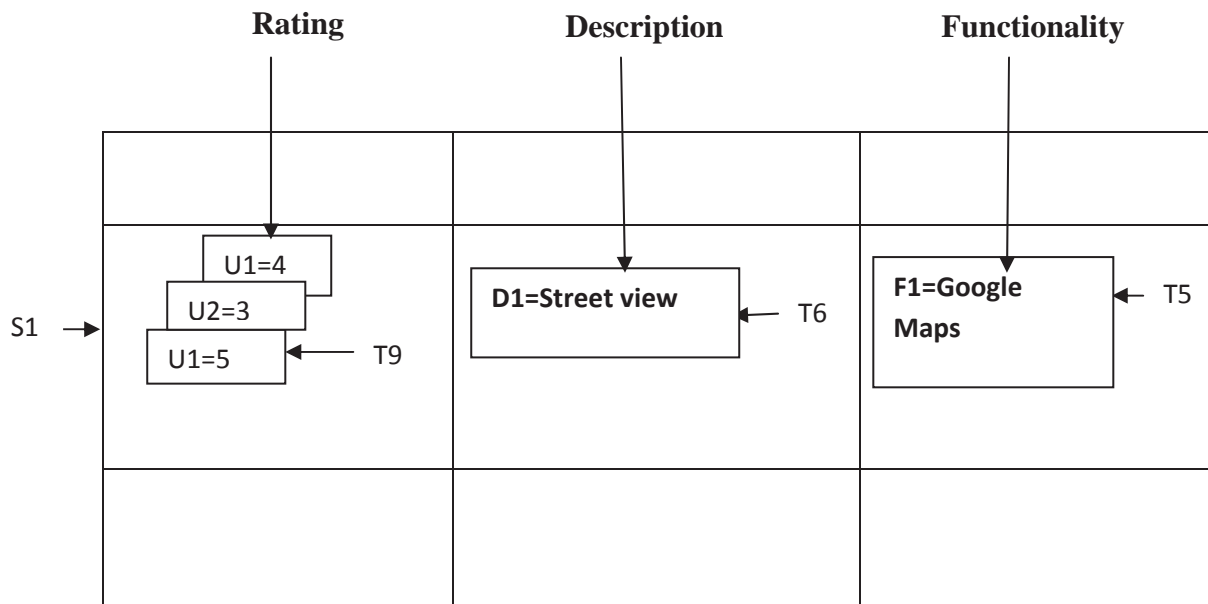


Table 1: Slice of Service Big Table

### III. PROPOSED SYSTEM

The proposed method is implemented using Clustering based Collaborative Filtering (ClubCF) approach, as it reduces the online computation time by clustering the services based on their similarities and then recommends the services to the active users. ClubCF approach focuses on two stages Clustering and the Collaborative Filtering.

Algorithm for ClubCF approach:

1. Calculate description similarity, functionality similarity and characteristic similarity between possible pairs of services.

1.1 Compute Description and Functionality similarities using ‘Jaccard similarity coefficient’ [4] by using Equations 1 and 2.

$$D\_sim(st, sj) = \frac{|D_i \cap D_j|}{|D_i \cup D_j|} \quad (1)$$

$$F\_sim(st, sj) = \frac{|F_i \cap F_j|}{|F_i \cup F_j|} \quad (2)$$

It can be inferred from the above Equations 1 and 2  $s_i, s_j$  referred to as services, consider an example of description similarity and functionality similarity between two mashup services 4 wheels route map ( $s_1$ ) and 100 destinations ( $s_3$ ) was

$$D\_sim(s_1, s_3) = \frac{|D_1 \cap D_3|}{|D_1 \cup D_3|} = 1/8 = 0.125$$

$$F\_sim(s_1, s_3) = \frac{|F_1 \cap F_3|}{|F_1 \cup F_3|} = 1/2 = 0.5$$

**1.2** Calculate characteristic similarity between  $s_i, s_j$  using Equation 3.

$$(s_i, s_j) = \alpha \times D\_sim(s_i, s_j) + \beta \times F\_sim(s_i, s_j) \quad (3)$$

In the above Equation,  $\alpha$  belongs to the weight of description similarity,  $\beta$  belongs to the weight of functionality similarity and  $\alpha + \beta = 1$ . The weights express the relative importance between two measures. Example of Characteristic similarity between two services  $s_1$  and  $s_3$  is

$$\begin{aligned} C\_sim(s_1, s_3) &= 0.5 \times D\_sim(s_1, s_3) + 0.5 \times F\_sim(s_1, s_3) \\ &= 0.5 \times 0.125 + 0.5 \times 0.5 = 0.3125 \end{aligned}$$

**1.3** Repeat steps 1.1 and 1.2 for all the pairs of services and form the distance matrix ( $D_m$ ).

**2** Apply Agglomerative hierarchical clustering algorithm [5], [6] on Distance matrix ( $D_m$ ) and group services with maximum distance as shown the procedure in below figure 1.

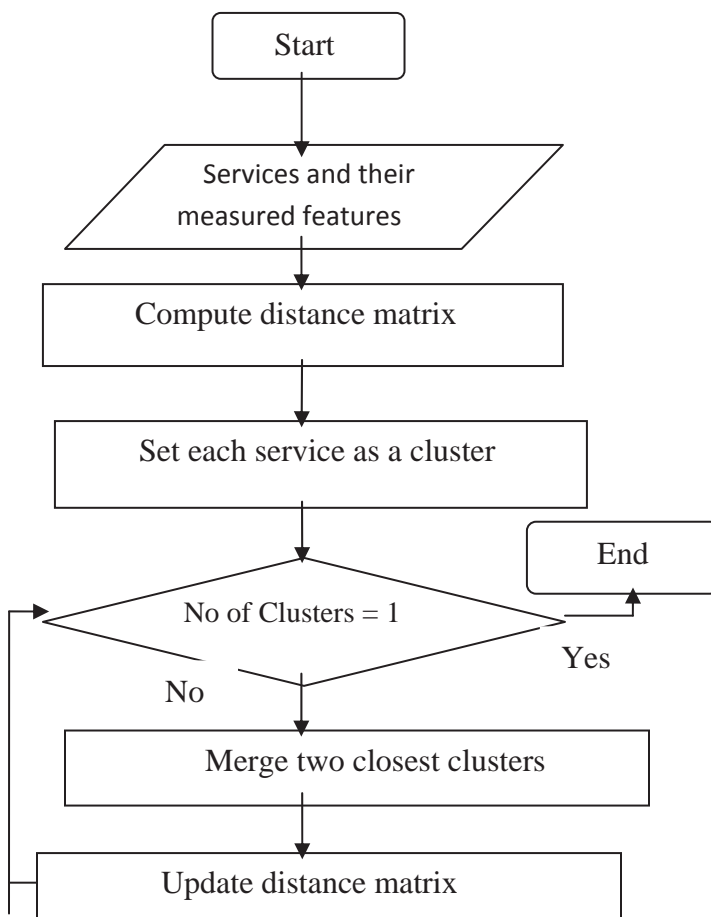


Fig 1: Flow Chart for AHC algorithm

3 Collect the Ratings given by the users to the service clusters formed in step 2.

3.1 Calculate rating similarities between the users using pearson correlation[7] coefficient given in Equation 4.

$$R_{sim}(s_i, s_j) = \frac{\sum_{u_i \in u_c \cap u_j} (r_{u_i, s_i} - \bar{r}_{s_i})(r_{u_i, s_j} - \bar{r}_{s_j})}{\sqrt{(\sum_{u_i \in u_c \cap u_j} (r_{u_i, s_i} - \bar{r}_{s_i})^2)(\sum_{u_i \in u_c \cap u_j} (r_{u_i, s_j} - \bar{r}_{s_j})^2)} \quad (4)$$

The Equation 4 depicts  $u_c$  is a set of users who rated  $s_i$  while  $u_j$  is set of users who rated  $s_j$ ,  $u_i$  is a user who rated both  $s_i$  and  $s_j$ ,  $r_{u_i, s_i}$  is the rating of  $s_i$  given by  $u_i$ .

3.2 If the number of users is limited then use enhanced rating similarities [8] for similarity calculations using Equation 5

$$R_{sim}(s_i, s_j) = \frac{2 * |u_c \cap u_j|}{|u_c| + |u_j|} * R_{sim}(s_i, s_j) \quad (5)$$

Where  $u_c \cap u_j$  is the number of users who rated both services  $s_i, s_j$ ,  $u_c$  and  $u_j$  are number of users who rated services  $s_i, s_j$ , respectively.

3.3 Calculate Neighbor's similarity of ratings for active users using the Equation 6.

$$N(s_i) = \{s_j | R_{sim}(s_i, s_j) > \gamma, s_i \neq s_j\} \quad (6)$$

The above Equation depicts  $R_{sim}(s_i, s_j)$  as the enhanced rating similarity between service  $s_i$  and  $s_j$ ,  $\gamma$  is the rating similarity threshold.

3.4 Compute predicted ratings for new users and also for the un rated users by using Equation 7.

$$r_{u_a, s_i} = \bar{r}_{s_i} + \frac{\sum_{s_j \in N(s_i)} (r_{u_a, s_j} - \bar{r}_{s_j}) * R_{sim}(s_i, s_j)}{\sum_{s_j \in N(s_i)} R_{sim}(s_i, s_j)} \quad (7)$$

Where  $\bar{r}_{s_i}$  is the average rating of  $s_i$ ,  $N(s_i)$  is

the neighbor set of  $s_i$ ,  $s_j \in N(s_i)$  denotes  $s_j$  is the neighbor of target service  $s_i$ ,  $r_{u_a, s_j}$  is the average rating of active user  $u_a$  gave to  $s_j$ ,  $\bar{r}_{s_j}$  is the average rating of  $s_j$ , and  $R_{sim}(s_i, s_j)$  is the enhanced rating similarity between services  $s_i, s_j$  computed using above Equation 7.

4 Calculate the MAE of ClubCF and compare the results with the IBCF approach.

#### IV. RESULTS

The proposed algorithm is implemented using Intel Dual Core processor with 2 GB RAM and Java (JDK 1.6) for coding. The experimental results of ClubCF approach are compared with the traditional Item based Collaborative Filtering which shows the computation time decreases as the number of services increases.

All the services are stored in the big table were labelled with row key column key, time stamp of their records along with their description and functionality. The major difference between the services storing in big table and the traditional Data Base Management systems was online computation time. The computation time decreases with in Big Table when compared to the DBMS storage mechanism. The following Figure shows the graphical representations of their comparison.

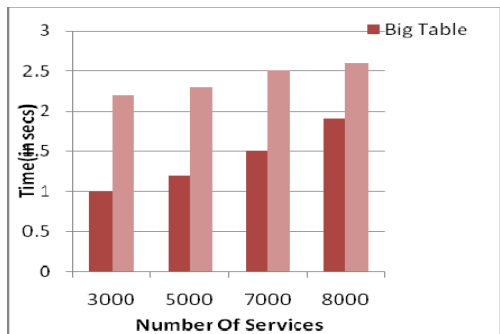


Fig 2: Big table vs traditional Data Base

The Big Table stores the mash up services along with their functionality and description whose similarities needs to be calculated in order to obtain the distance matrix which is the input for the Agglomerative hierarchical clustering algorithm.

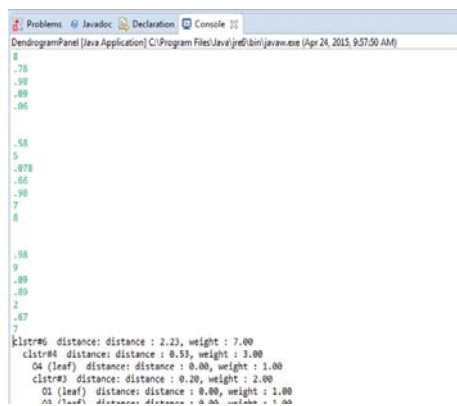


Fig 3: Input matrix

After providing the distance matrix as an input we will get the dendrogram panel representation of different clusters based on minimum or maximum distance. The following Figure represents the dendrogram panel of hierarchical clusters.

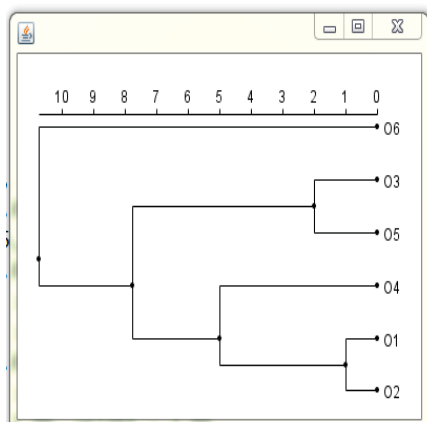


Fig 4: Dendrogram panel

After that the users need to rate the services with in the different clusters between 0 to 5 rating system. Then the recommendation mechanism came into picture in order to help the active users who don't know any thing about the service by calculating the rating similarities between different users and compute predicted ratings for new users.

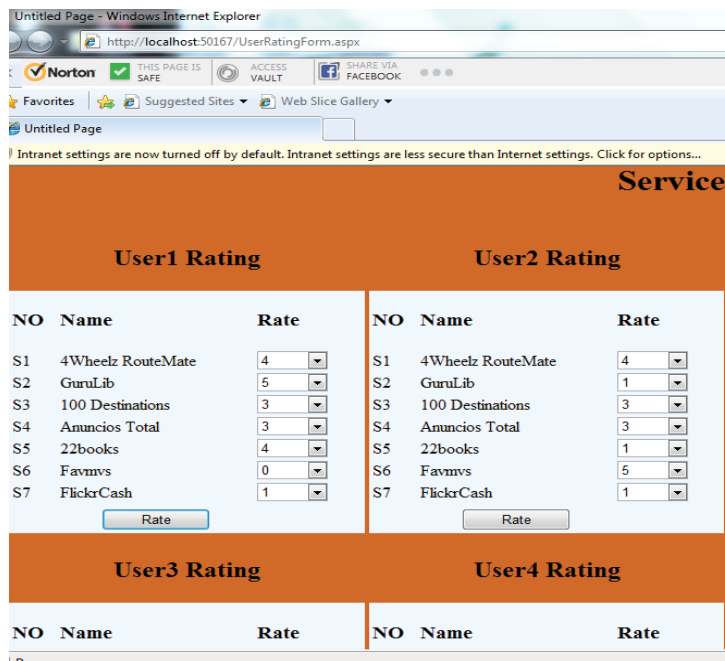


Fig 5: Rating for the services

All Ratings							
	S1	S2	S3	S4	S5	S6	S7
User1	4	5	3	3	4	0	1
User2	4	1	3	3	1	5	1
User3	3	4	1	0	2	1	3
User4	1	0	2	1	5	4	3

Fig 6: Ratings of users to services

```

1 package com.predictionmarketing.itemrecommend;
2
3 import java.io.File;
4
5
6
7
8
9
10 public class ItemRecommend {
11
12
13
14
15
16
17
18 public static void main(String[] args) {
19     try {
20         DataModel dm = new FileDataModel(new File("data/movies.csv"));
21         ItemSimilarity sim = new LogLikelihoodSimilarity(dm);
22         GenericItemBasedRecommender recommender = new GenericItemBasedRecommender(dm,sim);
23         int x=1;
24         for(LongPrimitiveIterator items = dm.getItemIDs();items.hasNext();){
25             long itemId = items.nextLong();
26         }
27     }
28 }

```

terminated-ItemRecommend [Java Application] C:\Program Files\Java\jre6\bin\javaw.exe (Apr 24, 2015, 10:08:29 AM)

```

3,100,0.98773706
4,56,0.99627966
4,174,0.99661305
4,204,0.9959589
4,202,0.99582237
4,385,0.9957967
5,216,0.99432045
5,308,0.9922024
5,234,0.99179345
5,56,0.99115413
5,53,0.9909523
6,547,0.97414243
6,923,0.97221303
6,221,0.97155356
6,1129,0.9711365
6.14.n.96879855

```

Fig 7: Recommendations

After calculating the rating similarities and generating recommendations for the un rated users we need to check for the accuracy of the ClubCF approach by comparing the mean absolute errors of ClubCF approach and the IBCF approach which leads the better accuracy of using the ClubCF approach. The below Figure represents the comparison between ClubCF and IBCF approaches.

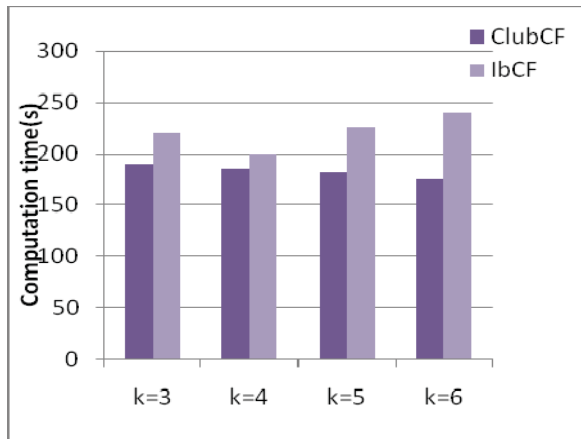


Fig 8: ClubCF vs IBCF approach

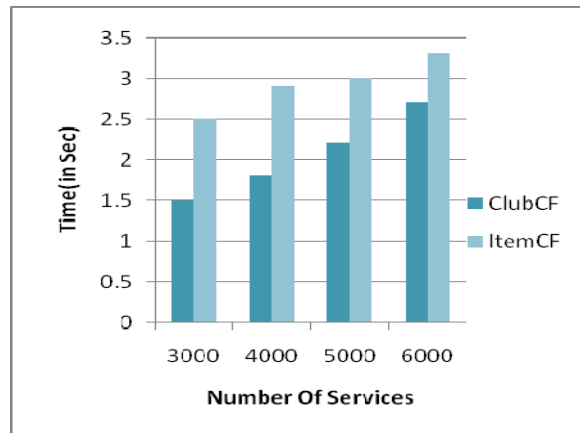


Fig 9: ClubCF vs IBCF Services

### V. CONCLUSIONS AND FUTURE WORK

The clustering of services is performed by using AHC algorithm over 6888 services. The actual rating similarities and predicted ratings are computed using PCC and enhanced rating similarities. The obtained results will be compared with the IBCF approach. The experimental results shows that computation time is approximately reduces by 30% in proposed method over IBCF even with increase in number of services. The proposed method also overcomes the problem of data Sparsity.

The future work can be extended to user based records and mining users interests and automatically provides recommendations to the users based on their explicit interests also, to develop methods to provide recommendations to new users even when few ratings are available.

## REFERENCES

- [1] Xiaoyuan Su and Taghi M. Khoshgoftaar .,“A survey of Collaborative Filtering Technique” A Review article, Hindawi Publishing Corporation Advances in Artificial Intelligence, Received 9 February 2009; Accepted 3 August 2009.
- [2] A. Yamashita, H. Kawamura, and K. Suzuki, “Adaptive Fusion Method for User-based and Item-based Collaborative Filtering,” *Advances in Complex Systems*, vol. 14, no. 2, pp. 133-149, May 2011.
- [3] F. Chang, J. Dean, S. mawat, et al.,“Big table: A distributed storage system for Structured data,” *ACM Trans. on Computer Systems*, vol. 26, no. 2, pp. 1-39, June2008.
- [4] A. Rodriguez, W. A. Chaovalitwongse, L. Zhe L, et al., “Master defect record retrieval using network-based feature association,” *IEEE Trans. on Systems, Man, and Cybernetics, Part C: Applications and Reviews*, vol. 40, no. 3, pp. 319-329, October 2010.
- [5] C. Platzer, F. Rosenberg, and S. Dustdar, “Web service clustering using multidimensional angles as proximity measures,” *ACM Trans. on Internet Technology*, vol. 9, no. 3, pp. 11:1-11:26, July, 2009.
- [6] G. Thilagavathi, D. Srivaishnavi, N. Aparna, et al., “A Survey on Efficient Hierarchical Algorithm used in Clustering,” *International Journal of Engineering*, vol. 2, no. 9, September 2013.
- [7] G. Adomavicius, and J. Zhang, “Stability of Recommendation Algorithms,” *ACM Trans. On Information Systems*, vol. 30, no. 4, pp. 23:1-23:31, August 2012.
- [8] D. Julie, and K. A. Kumar, “Optimal Web Service Selection Scheme With Dynamic QoS Property Assignment,” *International Journal of Advanced research in technology*, vol.2,no.2, pp. 69-75, may 2012.