**RESEARCH ARTICLE**                  **OPEN ACCESS**

# Analysis on Recommended System for Web Information Retrieval Using HMM

## Himangni Rathore[1], Hemant Verma[2]

[1] Computer Engineering Department, Vindhya Group of Technology and Science Khandwa Road, Indore, India
[2] Computer Engineering Department, Vindhya Group of Technology and Science Khandwa Road, Indore, India

*Abstract*

Web is a rich domain of data and knowledge, which is spread over the world in unstructured manner. The number of users is continuously access the information over the internet. Web mining is an application of data mining where web related data is extracted and manipulated for extracting knowledge. The data mining is used in the domain of web information mining is refers as web mining, that is further divided into three major domains web uses mining, web content mining and web structure mining. The proposed work is intended to work with web uses mining. The concept of web mining is to improve the user feedbacks and user navigation pattern discovery for a CRM system. Finally a new algorithm HMM is used for finding the pattern in data, which method promises to provide much accurate recommendation.

*Keywords*— Web mining, recommender systems, collaborative filtering, HMM.

## I. INTRODUCTION

Recommender system is an important part of information retrieval. They are used for enabling the users for filtering through large information and product data. It is based on information filtering system. An information filtering system is a system that avoids redundant or unwanted data from data which is streaming by user. Recommended systems commonly used in recent years and having several of applications. The most popular are music, researches, news, social sites, movies, jokes, online dating, financial insurance services etc. Sometimes recommender systems can apply techniques and methods from other areas such as information retrieval or human computer interaction. Recommender systems that constitute data mining techniques make their recommendations using knowledge learned from the actions and attributes of users. These systems are based on the user profiles that are based on two approaches item consumption or history data and the actions during the current session. The goal of these systems is to serve the right items to a user in a given context or create bundle package of same type of products to initialize the long term business objectives.

With the increasing of utility of internet accessing smart phones it is now possible to offer personalized, context-sensitive recommendations. This is difficult area of research as mobile data is more complex than the recommender system. It's promised to provide mobile users access to personalized recommendations anytime, anywhere.

## II. HISTORY

In 1985 for content filtering architecture a large scale information system is developed. In 1988 a rule based message filtering system has been proposed after that in 1990 an active mail-filtering agent called as MAFIA for an intelligent document processing support was developed by Lutz, E. R. In the history of recommendation system in 1992 Tapestry by Xerox Palo Alto came in exist. In 1994 the first system which is designed by collaborative filtering is Grouplens and after this in 1997 the first recommendation system using rating data is Movielens and it is the first movie recommender systems which provide a dataset for researchers. John S. Breese in 1998 analyse the empirical analysis of predictive algorithms in for the evaluation of user based collaborative filtering. Amazon proposed item based collaborative filtering in 2001 which patent is filed in 1998. Thomas Hofmann proposed PLSA in 1990 and applies same method on collaborative filtering in 2004 and Jonathan L. Herlocker also evaluates collaborative filtering recommender systems in the same year. After this the dynamic collaborative filtering has been proposed as ACM in Conference on Recommender System in 2007.

## III. BACKGROUND

Recommender system can be working in three ways- They are Collaborative filtering, content based filtering and Hybrid recommender system.

Collaborative filtering is an approach which is based on the collection and analyzing the large amount of information according to the users behavior, prediction, preferences and their activities what the users will like based on their similarities of

other users. There are mainly three approaches for collaborative filtering. They are memory based approach, model based approach and hybrid recommenders.

Memory based approach in collaborative filtering is an earlier mechanism and it is used in many business systems. It is very useful and effective method. Two memory based collaborative filtering are neighborhood based collaborative filtering and item based or user based collaborative filtering.

Model based collaborative filtering is developed to find the patterns by using different algorithms of data mining and machine learning. This collaborative filtering identifying the neighbors of an active user and prediction make according to the preferences on the new item. Some model based collaborative filtering are Bayesian networks, latent semantic models, clustering models.

Hybrid recommender system combines both the techniques memory based and model based collaborative filtering and improves the prediction performance.

| Categories of collaborative filtering | Techniques | Advantages | Disadvantages |
|---|---|---|---|
| Memory based CF | Neighbor based CF Item based top N recommendations. | No need to consider the content of item been recommended. Its implementation is easy. New data can be added easily. | It is dependent on human ratings. When data are sparse it decreases the performance. Limited scalability for large number of data sets. |
| Model based CF | Bayesian network, clustering model, latent semantic model, sparse factor analysis, probabilistic, markov decision process. | Prediction performance is improved. If the data are sparse it can be handling better. Having large amount of scalability for large sets. | Model building is expensive. Sometimes having confusion in scalability and prediction performance. By using reduction techniques sometimes it loses useful information. |
| Hybrid recommenders | Content based CF recommender. Hybrid CF combining memory based and model based CF algorithm. | It overcomes the limitation of CF content based or other recommendations. Prediction performance is improved. Overcome of scalability | Complexity and expensive have been increased for implementation. It usually needs external information which is not available easily. |

## IV. PROBLEM DOMAIN AND PROPOSED WORK

This section of the given document includes the problem domain and the solution that are proposed in order to optimize the available system.

As collaborative filtering includes two methods such as memory based or model based collaborative filtering. Many new methods are come to propose by theory or experiments. Despite the considerable research there is no clarity that which method is best. The performance of different methods is continuously differing from other based on the problem. Mainly performance in collaborative filtering is depending upon the number of users, number of items and sparsity level.

In the World Wide Web there are various different domain and subjects can be found in same place. In order to find the more appropriate data or subject in this complex domain is a frustrating process. The main reason behind this is exploration and investigation about each and every subject linearly. In place of exploring all data in a particular domain, suggestion about the interesting products,

subject and domain is provide ease in finding the appropriate data over any web space. Web recommender systems are basically designed to support users for finding the interest oriented products or make enable to find the similar products in huge data base for e-commerce applications. To design and develop a recommender system there are basically two things required first behavioral analysis and then a predictive model. Both of these methods are help to design an accurate and efficient recommender system. Therefore the recommendation system leads two main aspects first, personalization of web data according to the user need and second user behavior approximation. The proposed work is motivated by [3] concept where using rating and social score of rating is used to decide and recommend a product. This feedback of any product and rating this may helpful for finding the content in some of cases. But this system is not considered some of aspect of real time recommendation systems, some of them are personalization of the end client behavior and time based behavior fluctuation
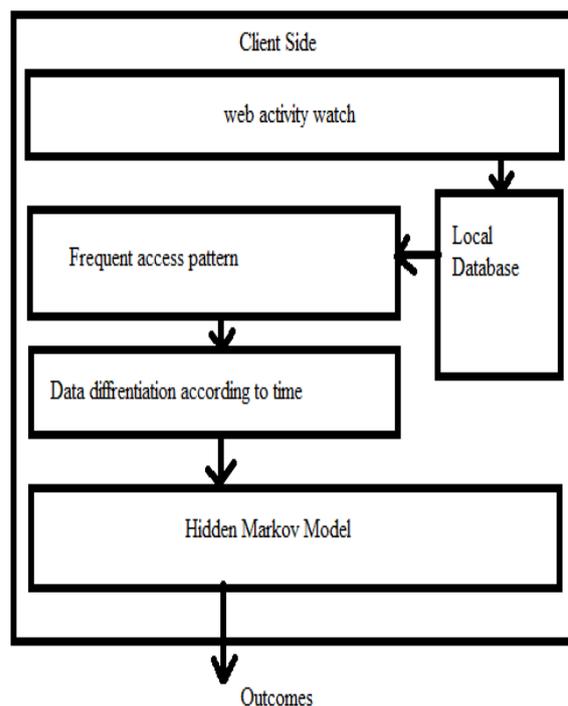
This section includes the primary objectives of the proposed study work, these task are required to accomplish during study.

**Study of different recommendation systems:** in this phase of study various research papers and articles are required to explore for recent development and analysis of new works in direction of the recommender system.

**Propose, design and develop a new recommender system for web search engine:** in this phase of study a web personalization and user behavior analysis a new recommender system is required to design, where for personalization collaborative filter is used for finding the frequent web access pattern of data and according to their data accessing pattern a hidden markov model is used for predicting the recommended object.

**Performance analysis and comparative study:** after implementation of the proposed methodology the performance evaluation of system is required for justification. The performance of recommender system is evaluated using the accuracy, and precision. On the other hand the required resources are also computed in order to get the cost (in terms of memory and time consumed) of recommender system. To overcome the above described problem a new solution is suggested for finding the optimum solution. The suggested method is given using the figure.

To simulate the proposed system a client end based solution is suggested in this work, where at client end web activity of the client system is evaluated first using the web activity watch. That may help to store the transactional patterns for personalization and behavioral analysis. Which is further stored in a local database for future analysis of data. A frequent access pattern algorithm is implemented with the local data base by which in different time slices the data is differentiated according to the time. The timely divided data is processed using hidden markov model in timely fashion for finding the pattern in data, which method promises to provide much accurate recommendation.



This section of the document provides the study work using which algorithm the proposed model becomes implementable.

An HMM is a double embedded stochastic process with two hierarchy levels. It can be used to model much more complicated stochastic processes as compared to a traditional Markov model. An HMM has a finite set of states governed by a set of transition probabilities. In a particular state, an outcome or observation can be generated according to an associated probability distribution. It is only the outcome and not the state that is visible to an external observer.

An HMM can be characterized by the following [7]:

1. N is the number of states in the model. We denote the set of states' $S = \{S1; S2;...SN\}$, where Si, i= 1;2;...;N is an individual state. The state at time instant t is denoted by qt.

2. M is the number of distinct observation symbols per state. The observation symbols correspond to the physical output of the system being modelled. We denote the set of symbols $V = \{V1; V2; ...VM\}$, where Vi, I = 1; 2; ... M is an individual symbol.

3. The state transition probability matrix A = [aij], where

$$a_{ij} = P\left(q_{t+1} = S_j \middle| q_t = S_i\right), 1 \leq i \leq N; t = 1,2$$

For the general case where any state j can be reached from any other state i in a single step, we have aij> 0 for all i, j. Also,

$$\sum_{j=1}^{N} a_{ij}, 1 \leq i \leq N$$

4. The observation symbol probability matrix B = {bj(k)}, where

$$b_j(k) = P(V_k|S_j), 1 \leq j \leq N, 1 \leq K \leq M$$

$$\sum_{k=1}^{M} b_j(k) = 1, 1 \leq j \leq N$$

5. The initial state probability vector$= \pi i$ , where

$$\pi_i = P(q_1 = S_i), 1 \leq i \leq N$$

such that

$$\sum_{i=1}^{N} \pi_i = 1$$

6. The observation sequence O = O1; O2; O3; ...OR, where each observation Ot is one of the symbols from V, and R is the number of observations in the sequence.

It is evident that a complete specification of an HMM requires the estimation of two model parameters, N and M, and three probability distributions A, B, and $\pi$. We use the notation $\wedge = (A; B; \pi )$ to indicate the complete set of parameters of the model, where A, B implicitly include N and M.

An observation sequence O, as mentioned above, can be generated by many possible state sequences. Consider one such particular sequence Q = q1; q2; ...; qR; where q1 is the initial state. The probability that O is generated from this state sequence is given by

$$P(O|Q, \lambda) = \prod_{t=1}^{R} P(O_t|q_t, \lambda)$$

Where statistical independence of observations is assumed Above Equation can be expanded as

$$P(O|Q, \lambda) = b_{q1}(O_1) . b_{q2}(O_2) ... b_{qR}(O_R)$$

The probability of the state sequence Q is given as

$$P(Q|\lambda) = \pi_{q1} a_{q1 q2} a_{q2 q3} \cdots a_{qR-1qQ}$$

Thus, the probability of generation of the observation sequence O by the HMM specified by can be written as follows:

$$P(O|\lambda) = \sum_{all\ Q} P(O|Q, \lambda)P(Q|\lambda)$$

Deriving the value of $P(O|\lambda)$ using the direct definition of is computationally intensive. Hence, a procedure named as Forward-Backward procedure is used to compute$P(O|\lambda)$.

The proposed work is evaluated and similarly various pre-existing models are also observed for finding the problem and solution which is appropriate for the given context of recommendation system design. In near future the proposed model is implemented using visual studio framework and the performance in terms of their predictive accuracy is measured in addition of that the resource consumption of the implemented system is also evaluated to estimate the efficiency of the given system.

## V. CONCLUSION

Web contains a lot of information and data, this data and knowledge can be used for extracting the information. Apart from the available data on web pages there is more information available in form of web access log and their structures. Therefore mining and extracting information from these logs and structures are also helpful for various different application system designs. The proposed work is based on the web uses mining data where system accepts the input according to the user behaviour and finding pattern to provide ease in search a specific kind of product. Implementation of the proposed solution for web based recommendation system. We are expecting an efficient and adoptable recommendation system which provides ease in product selection and recommendation for ecommerce application domains. In addition of that the proposed methodology is also helpful for the web search and page rank algorithms According to the time varying user behaviour fluctuations.

## VI. ACKNOWLEDGMENT

## REFERENCES

[1] *Extending Recommender Systems for Disjoint User/Item Sets: The Conference Recommendation Problem Mark F. Hornick*, Member, IEEE, and Pablo Tamayo, Member, IEEE Computer Society

[2] *Clustering Web Log Files – A Review*, R. Suguna, D. Sharmila, International Journal of Engineering Research & Technology (IJERT), Vol. 2 Issue 4, April – 2013

[3] *Study of Pre-processing Methods in Web Server Logs*, Dr. Sanjeev Dhawan, Mamta Lathwal, International Journal of Advanced

Research in Computer Science and Software Engineering, Volume 3, Issue 5, May2013

[4]  *Identifying User Behavior by Analyzing Web Server Access Log File*,K. R. Suneetha, Dr. R. Krishnamoorthi, IJCSNS International Journal of Computer Science and Network Security, VOL.9 No.4, April 2009

[5]  *Improving the Efficiency of Web Usage Mining Using K-Apriori and FP-Growth Algorithm*, Mrs. R. Kousalya, Ms. K. Suguna, Dr.V. Saravanan, International Journal of Scientific & Engineering Research Volume 4, Issue3, March-2013 ISSN 2229-5518

[6]  *Applying Web usage mining for personalizing hyperlinks in Web-based adaptive educational systems*, Cristóbal Romeroa, Sebastián Ventura, Amelia Zafra, Paul de Bra, 2009 Elsevier Ltd. All rights reserved

[7]  *Towards Understanding Learning Behavior Patterns in Social Adaptive Personalized E-Learning Systems*, Lei Shi, Alexandra I. Cristea, Malik Shahzad Awan, Craig Stewart, Maurice Hendrix, Proceedings of the Nineteenth Americas Conference on Information Systems, Chicago, Illinois, August 15-17, 2013.

[8]  *A survey on web recommender systems*, k. suneetha, p. sunil kumar reddy, publications of problems & application in engineering research – paper, vol 04, Special Issue01; 2013

[9]  *Study on Various Web Mining Functionalities using Web Log Files*, Supinder Singh, Sukhpreet Kaur, IJCSMC, Vol. 2, and Issue. 4, April 2013, pg.164 – 169

[10]  *Improved Approach to predict user Future Sessions using Classification and Clustering*, AkshayKansara, Swati Patel, International Journal of Science and Research (IJSR), India Online ISSN: 2319-7064