

An Efficient Web Recommender System based on Approach of Mining Frequent Sequential Pattern from Customized Web Log Preprocessing

Manisha Valera, GTU

Department of Computer Engineering,
C.U. Shah College Of Engineering and Technology,
Gujarat, India
Email: valeramanisha@gmail.com

Uttam Chauhan, GTU

Department Of Computer Engineering,
Vishwakarma Government Engineering College,
Gujarat, India
Email: ugchauhan@yahoo.com

Abstract—Internet has penetrated into every areas of society and has also become a huge, pervasive distribution and global information service center..In a real world The only option is to capture the attention of the user and provide them with the recommendation list to match the needs of user and keep their attention in their web site.. Web usage mining is a kind of data mining method that provide Smart personalized online services such as web recommendations, it is usually necessary to model users' web access behavior. Web usage mining includes three process, namely, preprocessing, pattern discovery and pattern analysis. The data reduction is achieved through data preprocessing. The aim of discovering frequent sequential access patterns in Web log data is to obtain information about the navigational behavior of the users. In the proposed system, an efficient sequential pattern mining algorithm is used to identify frequent sequential web access patterns. The access patterns are retrieved from a Graph, which is then used for matching and generating web links for recommendations.

Keywords—Web Mining, Sequential Pattern Mining, Recommendation, Web Usage Mining (WUM), Preprocessing, Pattern Discovery, Pattern Analysis, Weblog, Sequential Patterns.

I. INTRODUCTION

In this world of Information Technology, as days go , The Internet has affected almost every aspect of our world. Since the number of web sites and web pages has increased rapidly, discovering and understanding web users' surfing behavior are essential for the development of successful web monitoring and recommendation systems. During that , the substantial increase in the number of websites presents a challenging task for masters to manage the contents of websites to cater to the need of user's [1].

Web mining can be categorized into three areas of interest based on which part of the web to mine [3]:

1. Web Content Mining
2. Web Structure Mining
3. Web Usage Mining

Web usage mining is categorized into three phases:

A. Preprocessing :

According to the client, server and proxy server, the preprocessing is the first approach to retrieve the raw data from web resources and process the data. It automatically transforms the original raw data to the next process.

B. Pattern Discovery :

According to the data preprocessing, the raw data is used to discover the knowledge and to implement the techniques which will be used for machine learning. This makes use of data mining procedures.

C. Pattern Analysis :

It is the process after pattern discovery. It checks whether the pattern is correct on the web and guides the process of extraction of the information/ knowledge from the web.

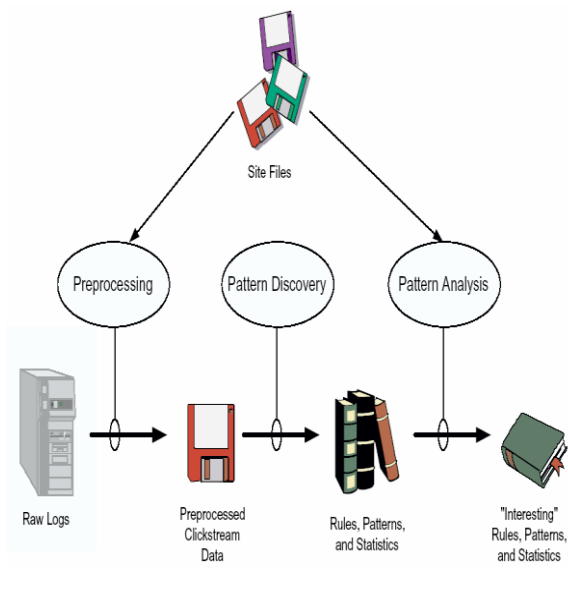


Fig. 1. Process Of Web Usage Mining

II. PROPOSED ALGORITHM

Sequential Web page Access pattern mining has been a focused theme in data mining research for over a decade with wide range of applications. The aim of

discovering frequent sequential access (usage) patterns in Web log data is to obtain information about the navigational behavior of the users. This can be used for advertising purposes, for creating dynamic user profiles etc

The proposed algorithm constructs a Graph to capture the user's web access behavior of a website and then uses the data mining steps in order to find out the Frequent Sequential Web Access Patterns. Then Web recommendation system is designed based on the Frequent Sequential Web Access Patterns. In this approach, firstly we are applying the pre-processing step on Input Web server logs to get the list of accessed web page sequence within different sessions, and then will construct a Web Usage Graph using accessed web page sequence in all sessions. Finally applying mining steps to mine the useful sequential user access patterns. From sequential user access patterns we generate recommendation of web page.

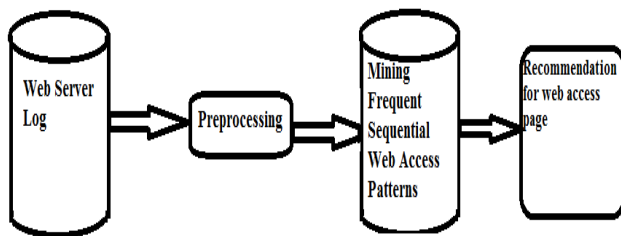


Fig. 2. Block diagram of the proposed web recommendation system

The purpose of data preprocessing is to extract useful data from web log and then transform these data in to the form necessary for pattern discovery. The input for the proposed web recommendation system is a web server log data and it comprises IP address, access time, HTTP request method used, URL of the referring page and browser name. It is difficult for these web server log data to be directly used to mine the desired sequential pattern mining process [10]. There are three steps in preprocessing of log data: Data cleaning, user identification, session identification.

The log entry contains various fields which need to be separate out for the processing. The process of separating field from the single line of the log file is known as field extraction. The server used different characters which work as separators. The most used separator character is or 'space' character.

The process of data cleaning is removal of outliers or irrelevant data. The Web Log file is in text format then it is required to convert the file in database format and then clean the file. First, all the fields which are not required are removed and finally we will have the fields like date, time, client ip, URL access, Referrer and Browser used/ Access log files consist of large amounts of HTTP server information. Analyzing, this information is very slow and inefficient without an initial cleaning task. All log entries with file name suffixes such as gif, JPEG, jpeg, GIF, jpg, JPG can be

eliminated since they are irrelevant [4]. Web robot (WR) (also called spider or bot) is a software tool that periodically a web site to extract its content[6]. To identify web robot requests the data cleaning module removes the records containing "Robots.txt" in the requested resource name (URL). The HTTP status code is then considered in the next process of cleaning by examining the status field of every record in the web access log, the records with status code except 200 are removed because the records with status code 200, gives successful response[7].

Using User Identification, we can identify individual user by using their IP address. If new IP address, there is new user. If IP address is same but browser version or operating system is different then it represents different user. User identification an important issue is how exactly the users have to be distinguished. It depends mainly on the task for the mining process is executed. In certain cases the users are identified only with their IP addresses .

The different sessions belonging to different users should be identified. Web logs can be regarded as a collection of sequences of access events from one user or session in timestamp ascending order. Here we are defining a time interval of 1 hour for each session [11].

Algorithm : Frequent Sequential Pattern Mining

Input: Web_Access_Sequence S – accessed page list of each session.

Min_count = No. Of Sessions * Percentage of support/100

Output: Frequent_Web_Access_Patern – Table of frequent_pattern with their respective frequency, arranged according to length of sequence.

Create_graph (Web_Access_Sequence S)

```

{
  For (Each Session Si in Web_Access_Sequence S)
  {
    For (Each webpage Pj in page sequence of Si)
    {
      If (i=j=1) {
        create_node (Pj)
      }
      Else If (Pj does not exists in all nodes generated so far)
      {
        create_node (Pj)
      }
      Else
      {
        set node_count(Pj) = node_count (Pj) + 1
      }
    }
    For (Each adjacent pair of webpage Pj in page sequence)
    {
      If (i=j=1) {
        create_edge (source node(Pj), destination node(Pj+1))
        set link_count= 1
      }
      Else If (an edge with same source, destination does not exists in all edges generated so far) {
        create_edge(source node(Pj), destination node(Pj+1))
        set link_count= 1
      }
    }
  }
}
  
```

```

}
Else If (an edge with same source, destination exists in all
edges generated so far in Si) {
    link_count = link_count + 1
}
}
}
}
prune_Graph(Web_Usage_graph <N,E>,min_count)
{
Remove_infrequent_nodes()
{
For (each session Si in web_access_sequence S)
{
    For(each node Pj in Si)
    {
        If(node_count(Pj)<min_count)
        {
            Find in_edge(Pj) and out_edge(Pj)
            create an edge E' with source node
            from in_edge(Pj) and destination from
            out_edge(Pj)
            If (an edge E with same source,
            destination exists as of E' in Prune_graph)
            {
                link_count(E) = link_count(E) + 1
            }
            else
            { link_count(E)=1}

            Remove node(Pj).
        }
    }
}
}
}
For (Each edge E in Pruned_Web_Usage_graph G)
{
    If(link_count(E) < min_count)
    remove edge E from Pruned_Web_Usage_graph G
}
}

Mining_graph(Pruned_Web_Usage_graph G')
{
For (All nodes & edges in Pruned_Web_Usage_graph G')
{
    Traverse the longest path
    set length=number of nodes in longest sequence
    while (length>0)
    {
        Traverse all paths with number of nodes=length,in
        Pruned_Web_Usage_graph G'
        add the visited nodes in the path as frequent_pattern
        add frequent_pattern to frequent_pages_list
        if(length==1)
        {
            set frequency of frequent_pattern=node_count
        }
        Else
        {
            set frequency of frequent_pattern=minimum of
            link_count of visited edges in the frequent_pattern
            length=length-1
        }
    }
    Return (frequent_pattern, frequency, length)
}
}

```

Algorithm: Web Page Recommendation Rules Generation

Input: Frequent Sequential Patterns

Output: RR – recommendation rule of a set of ordered access events for S.

1. frequent sequence with minimum support count.
S = a1a2... an - current access sequence of a user,
 2. *MinLength* - minimum length of access sequence,
MaxLength - maximum length of access sequence
 3. Initialize $RR = \emptyset$.
 4. If $|S| > \text{MaxLength}$ then remove the first $|S| - \text{MaxLength} + 1$ items from S.
 5. If $|S| < \text{MinLength}$ then return RR
 6. For each item a_i from S to the end do
 If Current item points to next item
 Then insert the next item into RR order by their Support.
 Else
 Remove the First item from S and repeat from Step5.
- Return RR.

In general, each line of web logs (one access record) includes the following key information: date-timestamp, client IP address, user ID, requested URL, and HTTP status code. Given a sequence database where each sequence is a list of transactions ordered by transaction time with each transaction comprising a set of items, find all sequential patterns with a user-specified minimum support, which is defined as the number of data sequences containing the pattern. Let E be a set of unique access events, which represents web resources accessed by users, i.e. web pages, URLs, topics or categories. A web access sequence $S = e_1e_2 \dots e_n$ ($e_i \in E$) for $1 \leq i \leq n$ is an ordered collection sequence of access events, and $|S| = n$ is called the length of the web access sequence. A web access sequence S is called a sequential web access pattern, if $\text{sup}(S) \geq \text{MinSup}$, where MinSup is a given support threshold. An access event $e_i \in E$ is called a frequent event, if $\text{sup}(e_i) \geq \text{MinSup}$. Otherwise, it is called an infrequent event. A Pattern-Data Structure model is proposed for storing sequential web access patterns compactly, so that it can be used for matching with a user's current access sequence and generating recommendation rules more efficiently in the Recommendation Rules Generation component.

Based on list of accessed web page sequence within different sessions, a Directed Graph can be constructed called web usage graph. Graph consists of vertices (nodes) and edges (links) in which, nodes are for web pages, and edges represent sequential access between the pages. The number of nodes in the graph is equal to number of distinct web pages accessed during all sessions, after creating a graph first we have to remove the infrequent node from the given web access sequence, which can be decided by checking that the occurrence of web page in each sequence as a page count is less than the given min_count then it is called infrequent. So we have to just remove the webpage from the given web access sequence. Each edge in

graph contains weight as link count, that represents the frequency of edge, edge id & list of sessions involved in that path or edge. In the $\text{prune_Graph}(\text{Web_Usage_graph}, G\langle N, E \rangle, \text{min_count})$, N is for nodes and E is for Edges. If the link count is less than the min_count then we have to remove the link. In the Process of Mining, By traversing nodes of Pruned Web Usage graph starting through each node, we get the frequently occurred sequence of nodes that represents frequent web access patterns. The frequency of the sequence will be the minimum link count of all the edges involved in that sequence. In such manner, traverse all the existing path in the Pruned Web Usage graph and enlist all frequent patterns along with their frequencies and arrange them by order of length. Frequent sequence with length of 1 are all the nodes in Pruned Web Usage graph itself, node count will represent their frequency. Following Algorithm is for mining the pruned web usage graph to get set of frequent sequential web access Patterns.

Recommender System is one kind of filtering system which is used for ranking or priority of the Webpage. The recommendations are retrieved for a given user's web access sequence S . length of the user web access sequence S must satisfy the thresholds (minlen and maxlen). If its length is greater than maxlength then we have to remove first ($S - \text{maxlength} + 1$) element. If it contains next item then return the recommendation rule order by the support [12].

III. IMPLEMENTATION AND RESULT ANALYSIS

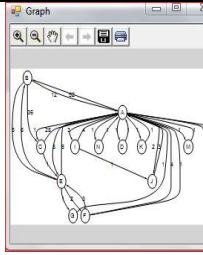
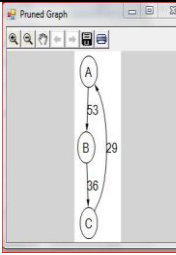
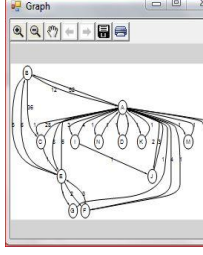
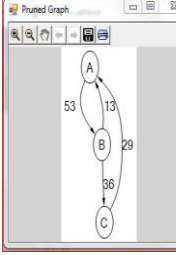
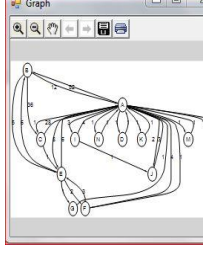
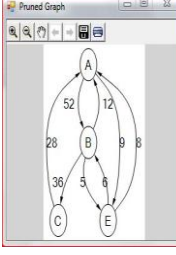
The log Files are Stored at IIS Server. For pattern discovery first we have to convert the Log File into SQL table .for that we have used JITBIT LOG2SQL software, which converts the log file into SQL table. We have used the Visual Studio .net framework 3.5 and SQL Server to preprocess the Web Log Data. We have taken the Web Log data of one site. The total number of records is 42944.

After that we have Performed Cleaning on the Data. So the Data is reduced. The Records' URL Stem which contains .jpg, .css, .js, .png, .gif, robot.txt, .bmp etc is removed.

Then we have done the User Identification. So, total distinct 23 users are found.

Then we have Performed Session Identification. So individual User Contains Session of one Hour. If his total hit time occurs more than 1 hour then new Session occurs. then, we store the session's web access sequence in a table and giving each page one alias name. for creating a web usage graph we are giving text file of web access sequence as an input. Then we create a graph. and after pruning it's structure will be reduced. so, complexity decreases. The pruned graph is applied as an input for the depth first search algorithm. From that we mine the frequent sequential web access patterns order by their length. From the frequent sequence patterns we get the recommendation rule order by their support.

TABLE I. TIME FOR GENERATING A GRAPH AND PRUNED GRAPH FROM OLD AND MODIFIED ALGORITHM

Min_count	Graph	Pruned Graph	Time to generate
15			Graph: 5.2499ms Pruned Graph: 12.4318ms
10			Graph: 5.04074ms Pruned Graph: 11.2491ms
5			Graph: 5.26273ms Pruned Graph: 10.9745ms

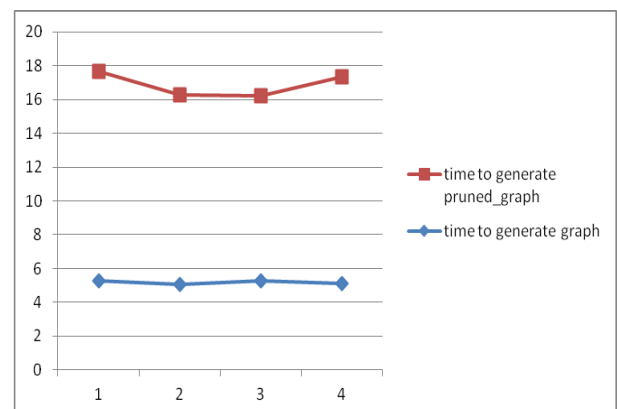


Fig. 3. time to generate graph and pruned graph

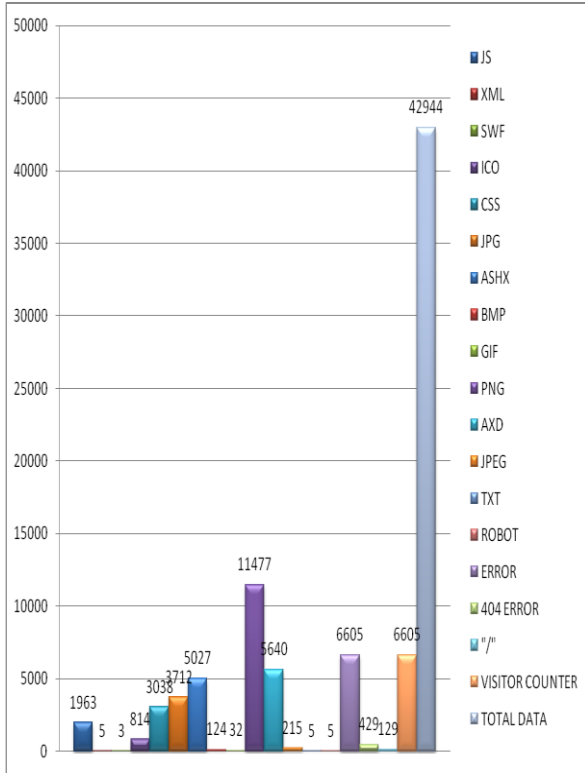


Fig. 4. Preprocessing of Web Log Data

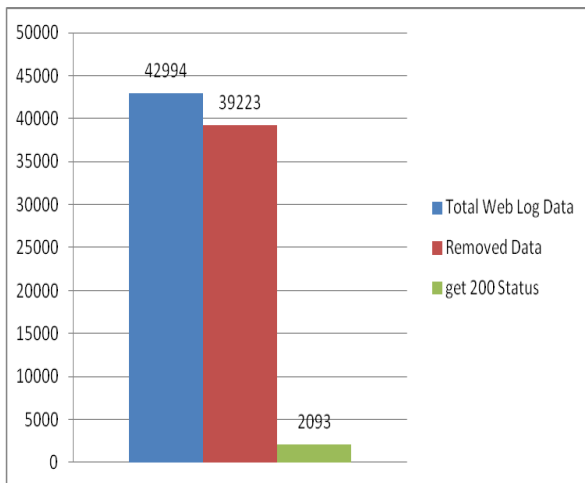


Fig. 5. after Preprocessing getting Customized Data for Use

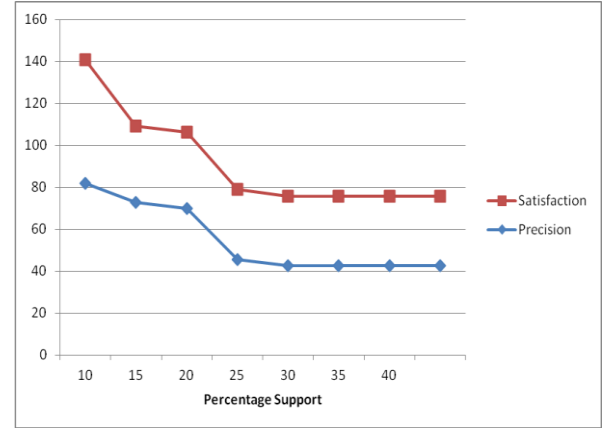


Fig. 6. Percentage support Vs precision and Satisfaction

A. Evaluation Measures

Let $S = a_1 a_2 \dots a_k a_{k+1} \dots a_n$ be a web access sequence. For the prefix sequence $S_{prefix} = a_1 a_2 \dots a_k$ ($k \geq \text{MinLength}$), we generate a recommendation rule $RR = \{e_1, e_2, \dots, e_m\}$ using the Graph, where all events are ordered by their support.

- If $a_{k+1} \in RR$, we label the recommendation rule as *correct*, and otherwise *incorrect*.
- If there exists $a_i \in RR$ ($k+1 \leq i \leq k+1+m$, $m > 0$), we deem the recommendation rule as *m-step satisfactory*. Otherwise, we label it as *m-step unsatisfactory*.

Let $R = \{RR_1, RR_2, \dots, RR_n\}$ be a set of recommendation rules, where RR_i ($1 \leq i \leq n$) is a recommendation rule. $|R| = n$ is the total number of recommendation rules in R . We define the following measures.

Definition 1: Let R_c be the subset of R that consists of all correct recommendation rules. The overall web Recommendation *precision* is defined as $\text{Precision} = |R_c| / |R|$.

This precision measures how probable a user will access one of the recommended pages.

Definition 2: Let $R_s(m)$ be the subset of R that consists of all *m-step* satisfactory recommendation rules. The *m-step Satisfaction* for web recommendation is defined as $\text{Satisfaction} = |R_s(m)| / |R|$.

In Preprocessing we are cleaning the data. So, the data is reduced, it takes less space of memory storage.

The experimental results given in Fig.10 show that as the Percentage Support threshold is increased, the *precision* and *satisfaction* remain relatively stable. From both experiments, we can conclude that better recommendations can be obtained with smaller support thresholds, at the expense of increased computational complexity for sequential web access pattern mining and maintaining a larger.

IV. CONCLUSION

Recommendation for a user can be achieved through web usage mining. Common navigation behaviours of the users can be used to improve the actual design of web pages and for making other modifications to a Web site. The contribution of the paper is to introduce a new way of web usage mining, and to show how frequent pattern discovery tasks can be accomplished by capturing complex user's browsing

behaviour into a graph in order to obtain hidden useful user's access patterns information.

REFERENCES

- [1] S.K. Pani, L.Panigrahy, V.H.Sankar, Bikram Keshari Ratha, A.K.Mandal, S.K.Padhi, "Web Usage Mining: A Survey on Pattern Extraction from Web Logs", International Journal of Instrumentation, Control & Automation (IJICA), Volume 1, Issue 1, 2011
- [2] Yogish H K, Dr. G T Raju, Manjunath T N, "The Descriptive Study of Knowledge Discovery from Web Usage Mining", IJCSI International Journal of Computer Science Issues, Vol. 8, Issue 5, No 1, September 2011
- [3] Udayasri.B, Sushmitha.N, Padmavathi.S, "A LimeLight on the Emerging Trends of Web Mining", Special Issue of International Journal of Computer Science & Informatics (IJCSI), ISSN (PRINT):2231-5292,Vol.-II,Issue-1,2
- [4] Navin Kumar Tyagi, A.K. Solanki & Sanjay Tyagi. "An Algorithmic approach to data preprocessing in Web usage mining", International Journal of Information Technology and Knowledge Management July-December 2010, Volume 2, No. 2, pp. 279-283
- [5] Surbhi Anand , Rinkle Rani Aggarwal, "An Efficient Algorithm for Data Cleaning of Log File using File Extensions", International Journal of Computer Applications (0975 – 888),Volume 48– No.8, June 2012
- [6] J. Vellingiri and S. Chenthur Pandian, "A Novel Technique for Web Log mining with Better Data Cleaning and Transaction Identification", Journal of Computer Science7 (5): 683-689,2011 ISSN 1549-3636 © 2011 Science Publications
- [7] Priyanka Patil and Ujwala Patil, "Preprocessing of web server log file for web mining", National Conference on Emerging Trends in Computer Technology (NCETCT-2012)", April 21, 2012
- [8] Vijayashri Losarwar, Dr. Madhuri Joshi, "Data Preprocessing in Web Usage Mining", International Conference on Artificial Intelligence and Embedded Systems (ICAIES'2012) July 15-16, 2012 Singapore
- [9] Sachin yele, Beerendra Kumar, Nitin Namdev, Devilal Birla, Kamlesh Patidar., "Web Usage Mining for Pattern Discovery", International Journal of Advanced Engineering & Applications, January 2011.
- [10] Chetna Chand, Amit Thakkar, Amit Ganatra, "Sequential Pattern Mining: Survey and Current Research Challenges", International Journal of Soft Computing and Engineering (IJSCE) ISSN: 2231-2307, Volume-2, Issue-1, March 2012
- [11] Dheeraj Kumar Singh, Varsha Sharma, Sanjeev Sharma "Graph based Approach for Mining Frequent Sequential Access Patterns of Web pages", International Journal of Computer Applications (0975 – 8887) Volume 40– No.10, February 2012.
- [12] Zhou.B, Hui. S.C and Chang. K, "An intelligent recommender system using sequential Web access patterns", IEEE Conference on Cybernetics and Intelligent Systems, vol. 1, pp. 393 - 398, December 2004.