

Knowledge Discovery from Web Logs -A Survey

S. Chitra¹ and Dr. B. Kalpana²

¹ Computer Science Dept, Avinashilingam University, Coimbatore, Tamilnadu

Received: 18 August 2011 Accepted: 7 September 2011 Published: 19 September 2011

Abstract

Web usage mining is obtaining the interesting and constructive knowledge and implicit information from activities related to the WWW. Web servers trace and gather information about user interactions every time the user requests for particular resources. Evaluating the Web access logs would assist in predicting the user behavior and also assists in formulating the web structure. Based on the applications point of view, information extracted from the Web usage patterns possibly directly applied to competently manage activities related to e-business, e-services, e-education, on-line communities and so on. On the other hand, since the size and density of the data grows rapidly, the information provided by existing Web log file analysis tools may possibly provide insufficient information and hence more intelligent mining techniques are needed. There are several approaches previously available for web usage mining. The approaches available in the literature have their own merits and demerits. This paper focuses on the study and analysis of various existing web usage mining techniques.

Index terms— Web Usage Mining, Personalization, Pre -processing, Web Log, Navigation Patterns.

1 INTRODUCTION

HE World Wide Web (WWW), is the current era of information explosion, has become the large source of online data, which includes text, graphics, videos, sound, etc. WWW is a comprehensive information medium in which the users can read, write and communicate through the use of computers connected to the Internet. Recent studies have estimated that the Web have more than one billion pages. It has become the powerful technology for sharing new ideas and content exchange. The impact of the Internet on everyday life is tremendous and it has changed the way of doing business, providing and receiving education, organization management, etc. The manner of information collection and sharing has changed with the advancement of hardware and communication software.

The growth has motivated the web service providers to predict the user's web usage behaviors so that, they can ? Personalize the information provided to them ? Make the websites more user friendly ? Reduce the traffic load ? Create or modify their website to suit different group of people.

The current requirement focuses on some tools which will help system analysts and business persons to learn user/consumer's needs, so that user requirements or demand can be solved immediately.

Web mining is the application of data mining techniques to web-based data for the purpose of learning or extracting knowledge. The techniques in web mining focus on providing solutions to content provider, web designer and programmers to improve their website and also to the web users with navigation assistance tools. It is a part of data mining where knowledge is gained from WWW.

Web servers trace and gather information about user interactions every time the user requests for particular resources. Evaluating the Web access logs would assist in predicting the user behavior and also assists in formulating the web structure.

Web usage mining, popularly also known as web log mining, works on the secondary data like web log file, click streams to extract knowledge with regard to web usage. It is the process which uses data mining techniques

2 LITERATURE SURVEY

44 abundantly, the result of which can be used for several uses like personalization, system improvement and site
45 modification.

46 It is essential to investigate what kind of features a WUM system is estimated to have with the intention of
47 performing effective and efficient Web usage mining, and what kind of challenges may be faced in the process of
48 developing new Web usage mining techniques. A Web usage mining system should be able to: ?

49 2 LITERATURE SURVEY

50 Web usage mining is a major application of data mining technology to extract the data of the Web server log
51 file. It can find out the browsing behavior of user and a certain type of correlations among the web pages. Web
52 usage mining offers the support for the Web site design, and also offers personalization server and further business
53 decision making, etc. Web mining applies the data mining, the artificial intelligence and the chart technology and
54 so on to the Web data and tracks the users browsing characteristics, and then obtains the users using pattern.
55 Qingtian Han et al., [1] investigated on Web Mining Algorithm based on Usage Mining. And it also provides the
56 design approach of the electronic commerce website application algorithm. This approach is easy, efficient and
57 easy to understand, it is appropriate to the Web usage mining demand of building a low cost B2C website.

58 In order to enhance the Web site, it is necessary to estimate current usage of the particular website. Web
59 usage mining and statistical analysis are two approaches to estimate usage of Web site. The integration of Web
60 usage mining and statistical analysis provides more exact information about Web usage. With the help of Web
61 usage mining approaches, graph mining focuses complex Web browsing patterns like parallel browsing. With the
62 use of statistical analysis approaches, investigating the page browsing time provides precious information about
63 Web site and its user's behavior. Heydari et al., [2] presented a Web usage mining technique which integrates
64 Web usage mining and statistical analysis by taking into consideration of client side data. In additional way, it
65 integrates graph based Web usage mining and browsing time examination with the consideration of client side
66 data. It assists in rebuilding user session precisely as it has been and in accordance with these data, this Web
67 usage patterns offers better accuracy.

68 Web usage mining has turn out to be very popular in different business fields associated with Web site
69 development. In Web usage mining, frequently browsed navigational paths are obtained with the assistance of
70 Web page addresses from the Web server visit logs, and the patterns are utilized in different applications together
71 with recommendation. Normally the semantic information of the Web page contents is not considered in Web
72 usage mining. Salin et al., [3] provided a structure for combining semantic information with Web usage mining.
73 The common navigational patterns are obtained as the form of ontology instances rather than Web page addresses
74 and the outcome is used for generating Web page recommendations to the visitor. Additionally, an assessment
75 method is implemented with the intention of testing the accomplishment of the recommendation. Test result
76 confirms that precise recommendations can be achieved by including semantic information in the Web usage
77 mining.

78 Nasraoui at al., [4] presented a comprehensive structure and findings in mining Web usage patterns by using
79 Web log files of a real Web site that has all the demanding aspects of real-life Web usage mining, together
80 with evolving user profiles and external data describing an ontology of the Web content. Therefore, the authors
81 present a technique for determining and tracing the mounting user profiles. The authors also discuss how the
82 obtained users profiles can be improved with clear information obtained from search queries of Web log data.
83 Profiles are also enhanced with additional domain-specific information aspects that provide a panoramic view of
84 the discovered mass usage modes. Many experiments have been done by the author to assess the excellence of
85 the mined profiles, especially their adaptability in the face of developing user behavior.

86 Web usage mining utilizes data mining methods to examine the user access of Web sites. As with any KDD
87 (knowledge discovery and data mining) process, WUM comprises of three main phases: preprocessing, knowledge
88 extraction, and results analysis. Tanasa at al., [5] concentrates on data preprocessing, a difficult and complicated
89 phase. Analysts intend to find out the accurate list of users who browsed the Web site and to reconstitute
90 user sessions-the order of actions every user carried out on the Web site. Inter-sites WUM focuses on the Web
91 server logs from numerous Web sites, commonly belonging to the similar organization. Therefore, analysts must
92 reconstruct the users' path through all the different Web servers that they browsed. This solution is to integrate
93 all the log files and reconstitute the visit. Traditional data preprocessing comprises of three phases: data fusion,
94 data cleaning, and data structurization. The author calls this solution as advanced data preprocessing. This
95 technique comprises of a data summarization phase, which will permit the analyst to choose only the information
96 importance. The authors have effectively tested this technique in an experimentation with log files from INRIA
97 Web sites.

98 Jianxi et al., [6] provides a Web usage mining method based on fuzzy clustering in identifying target group.
99 Data mining is a procedure of non-trivial mining of inherent, formerly unidentified, and extremely useful data
100 from very large quantity of data. Web mining can be defined mainly as the utilization of data mining techniques
101 to Web data. Web usage mining is a noteworthy and fast developing area of Web mining where several researches
102 has been carried our earlier.

103 The author utilized the fuzzy clustering method for identifying groups that allocate comparable interests and
104 behaviors by investigating the data gathered in Web servers.

105 Internet and Web technologies are extensively available, enabling it simpler for organizations to carry out

106 business and transfer data to customers. Furthermore, they accelerate financial transactions competently by
107 decreasing the transaction costs of commercial actions that businesses would generally incur. As a result, Internet
108 business has generated aggressive surroundings, a flourishing organizations wanted to survive and increase a
109 competitive advantage must offer a satisfactory package of customized services that convince customers' needs.
110 Regardless of the Internet's apparent benefits as a novel communication medium its advertising provides the same
111 advertising information to all customers and so has experienced from poor reactions. To increase a Web ad's
112 usefulness, Sung Min Bae et al., [7] developed a Web ad selector with the intention of personalizing advertising
113 information for customers according to their preferences and interests. The Web ad selection method segregates
114 the Web site customers with comparable preferences into numerous segments through Web usage mining. It makes
115 use of fuzzy rules that conveys customer segments' surfing patterns on the basis of specialist recommendation,
116 and recommends suitable ads by fuzzy inference.

117 Wu et al., [8] recommended a Web Usage Mining technique according to the sequences of clicking patterns in
118 a grid computing environment. Predicting user's browsing behavior is an important process of web usage mining.
119 It can support the web designers to improve the web structure or enhance the performance of the web servers.
120 Mining on the sequences of such clicking patterns (MSCP) can be considered as a data mining operation. MSCP
121 is normally an expensive process because of its considerable quantity of time for computation and storage for
122 archiving a huge quantity of information. Executing MSCP turns out to be unsuccessful or even not realistic on
123 a computer with limited resources. The author finds out the handling of MSCP in a distributed grid computing
124 environment and expresses its efficiency by experimental cases.

125 Web usage mining is a part of data mining technology to extract the data of the Web server log files. It can
126 find out the session patterns of user and certain kinds of correlations among these Web pages. Web usage mining
127 offers the support for the Web site design, by offering personalization server and additional business making
128 decision. There are several session patterns accumulated in Web server log files, page attribute of the same is
129 in Boolean quantity. With the purpose of enhancing the effectiveness of presented algorithms and decrease the
130 time of scanning database, and consequently focusing to these aspects, Gang Fang at al., [9] proposes a double
131 algorithm of Web usage mining in accordance with the sequence number that is appropriate for mining several
132 session patterns. The algorithm transforms session pattern of particular user into binary, and subsequently uses
133 up and down search approach to double generate candidate frequent itemsets. The algorithm works out support
134 by sequence number dimension with the intention of scanning once session pattern of a particular user, which is
135 dissimilar from conventional double search mining algorithm. in addition to this, the effectiveness of Web usage
136 mining is competently enhanced because of this approach. The experiment result confirms that the efficiency is
137 more rapid and more competent than the similar algorithms.

138 A lot of models are available and practices that analyze user behavior according to their user navigation
139 data and use clustering algorithms to differentiate their access patterns. The navigation patterns recognized are
140 predicted to satisfy the user's interests. Raghavendra et al., [10] modeled user behavior as a vector of the time
141 the particular user spends at each URL, and additionally categorize a new user access pattern. The clustering
142 and classification methods of k-means with non-Euclidean similarity measure, and artificial neural networks with
143 consistent inputs were implemented and evaluated. Despite recognizing user behavior, this model can also be
144 utilized as a prediction system in which it can be used to identify deviational behavior.

145 In Web Usage Mining (WUM), web session clustering plays a significant key part to categorize web visitors
146 according to the user click history and comparison measure. Swarm dependent web session clustering assists in
147 several ways to handle the web resources efficiently such as web personalization, schema modification, website
148 modification and web server performance. Hussain et al., [11] propose a structure for web session clustering at
149 initial level of web usage mining. The structure will cover the data preprocessing phase to organize the web
150 log data and transform the categorical web log data into numerical data. A session vector is acquired, so that
151 suitable comparison and swarm optimization possibly will be © 2011 Global Journals Inc. (US) applied to cluster
152 the web log data. The hierarchical cluster based technique will improve the existing web session techniques for
153 additional structured information about the user sessions.

154 With the huge development of World Wide Web and e-commerce the investigation of users' navigation patterns
155 has developed into more significant. Predicting users' navigation behavior is a challenging subject for ecommerce
156 enterprises. Web usage mining approaches can be utilized for modeling and predicting users' navigation patterns.
157 In actual fact, mining users' navigation pattern is the basic approach for producing recommendations. In practice,
158 the user interests are unpredictable, and it is complicated to follow the exact user navigation pattern. Khosravi
159 et al., [12] proposed a technique based on naive Bayesian method for modeling and predicting users' navigation
160 behavior. The author has used Web server logs as source data, and carried out his experiment.

161 Huge volumes of data are collected automatically by Web servers and accumulated in access log files.
162 Examination of server access data can offer important and helpful information. Web Usage Mining is the
163 method of using data mining approaches to find the usage patterns by using the Web data and is aimed towards
164 applications. It extracts the secondary data based on the interactions of the users throughout certain amount
165 of Web sessions. Web usage mining contains three stages, that is, preprocessing, pattern discovery, and pattern
166 analysis. Web usage mining has seen a huge increase of interest from the research people together with practice
167 communities. Etminani et al., [13] applied Kohonen's SOM (Self Organizing Map) to pre-processed Web logs of

4 CONCLUSION

168 Web server logs and extract frequent patterns. Experimental result of this technique confirms that this technique
169 would be more useful for Web site owner.

170 In order to offer the online prediction efficiently, Shinde et al., [14] formulated a architecture for online
171 recommendation for predicting in Web Usage Mining System. This approach provides the structural design
172 of on-line recommendation system in Web usage mining (OLRWMS) for enhancing the exactness of classification
173 by dealing between classifications, estimation, and provides user activities and user profile in online phase of this
174 architecture.

175 Nowadays, Internet has turned out to be an essential tool for each individual, in the same way Web usage
176 mining becomes a hotspot, which uses huge amounts of data in the Web server log and other appropriate data
177 sets for mining analysis and achieves valuable knowledge model about usage of relevant Web site. At the moment,
178 numerous works have to be performed with the positive association rules in Web usage mining, other than negative
179 association rules is considerably more significant, Yang Bin at al., [15] have applied negative association rules
180 approach to Web usage mining, in the course of the research the author have proved that the negative association
181 rules have an additional significant role on access pattern to Web visitors, provide the mining algorithms, to
182 resolve the deficiencies in which positive association rules are referred.

183 The recent development in Web technology has mounted the users and web pages at an exponential speed. The
184 evolutionary modifications in tools have made it promising to confine the users' concentrate and communications
185 with web applications through web server log file. Web log data is stored as text (.txt) file. Because of huge
186 quantity amount of unrelated information in the web log, the initial log data can not be straightforwardly
187 utilized in the web usage mining (WUM) process. As a result the preprocessing of web log data turns out to
188 be very important. The appropriate examination of web log file is valuable to control the web sites efficiently
189 for organizational and users' potential. Web log preprocessing is preliminary required process to enhance the
190 quality and efficiency of the future processing of web usage mining. There are number of methods existing at
191 preprocessing level of web usage mining. Various methods are utilized at preprocessing level like data cleaning,
192 data filtering, and data integration. Hussain at al., [16] analyzed the existing the preprocessing methods to
193 recognize the concerns and how WUM preprocessing can be enhanced for pattern mining and examination.

194 III.

3 PROBLEMS AND DIRECTIONS

195 The main objective of web usage mining is to recognize the interesting web usage patterns. In order to recognize
196 the interesting web usage patterns, a lot of researches are needed. They researches may focus on the following:

198 ? To provide precise page recommendation, it is necessary to understand the browsing behavior of the user
199 and it can be effectively done using the Machine Learning algorithms. This would definitely help in providing
200 the accurate page recommendations according to the user needs.

201 ? Based on the statistical analysis of the user, specifically, the amount of time user spending on a particular
202 page, the kind of links on which the user is interested, number of visit on a particular page will help to understand
203 the behavior of an user. ? The clustering technique can be used in grouping the user access pattern plays a
204 significant role in determining the resemblance in used browsing sessions. Therefore, the clustering technique can
205 be improved to enhance the performance of grouping the user sessions.

206 © 2011 Global Journals Inc. (US)

207 IV.

4 CONCLUSION

208 Web log files play a significant role in the Web Usage Mining. An important knowledge that can be obtained
209 from web log files is the user's navigation pattern. The challenge in obtaining such knowledge is that the users
210 are constantly shifting their focus and different users have different navigational behavior with different needs
211 associated with them. The navigation pattern knowledge can be used to help users by predicting their future
212 request and it will help on the personalization of websites. This paper provides the need for Web Usage Mining
213 and various techniques which focus on the Web Usage Mining. The directions provided in this paper will assist
214 the researchers to perform research on Web Usage Mining.^{1 2 3}
215

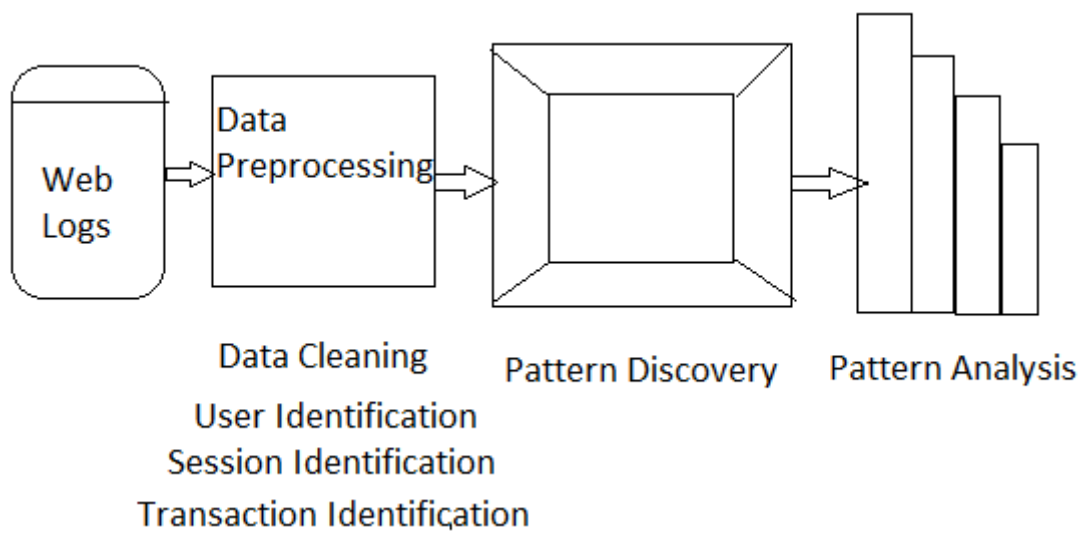
¹© 2011 Global Journals Inc. (US)

²OctoberKnowledge Discovery from Web Logs -A Survey

³OctoberKnowledge Discovery from Web Logs -A Survey



Figure 1:



1

Figure 2: Figure 1 :

-
- 216 [Bae] , Sung Min Bae .
217 [Bin] , Yang Bin .
- 218 [Fang and Xiong ()] ‘A Double Algorithm of Web Usage Mining Based on Sequence Number’. Gang Fang , ;
219 Jia-Le Wang; Hong Ying; Jiang Xiong . *International Conference on Information Engineering and Computer*
220 *Science (ICIECS)*, 2009. p. .
- 221 [Heydari et al. ()] ‘A graphbased web usage mining method considering client side data’. M Heydari , R A Helal
222 , K I Ghauth . *International Conference on Electrical Engineering and Informatics (ICEEI '09)*, 2009. 1 p. .
- 223 [Hussain et al. ()] ‘A hierarchical cluster based preprocessing methodology for Web Usage Mining’. T Hussain ,
224 S Asghar , S Fong . *6th International Conference on Advanced Information Management and Service (IMS)*,
225 2010. p. .
- 226 [Shinde and Kulkarni ()] ‘A New Approach for on Line Recommender System in Web Usage Mining’. S K Shinde
227 , U V Kulkarni . *International Conference on Advanced Computer Theory and Engineering*, 2008. p. .
- 228 [Nasraoui et al. ()] ‘A Web Usage Mining Framework for Mining Evolving User Profiles in Dynamic Web Sites’.
229 O Nasraoui , M Soliman , E Saka , A Badia , R Germain . *IEEE Transactions on Knowledge and Data*
230 *Engineering* 2008. 20 (2) p. .
- 231 [Tanasa and Trousse ()] ‘Advanced data preprocessing for intersites Web usage mining’. D Tanasa , B Trousse .
232 *IEEE Intelligent Systems* 2004. 19 (2) p. .
- 233 [Raghavendra et al. ()] ‘Comparative study of neural networks and kmeans classification in web usage mining’.
234 P S Raghavendra , S R Chowdhury , S V Kameswari . *International Conference for Internet Technology and*
235 *Secured Transactions (ICITST)*, 2010. p. .
- 236 [Khosravi and Tarokh ()] ‘Dynamic mining of users interest navigation patterns using naive Bayesian method’.
237 M Khosravi , M J Tarokh . *IEEE International Conference on Intelligent Computer Communication and*
238 *Processing*, 2010. p. .
- 239 [Ha ()] *Fuzzy Web ad selector based on Web usage mining*, Sang Chan Park; Sung Ho Ha . 2003. 18 p. .
- 240 [Han; Xiaoyan Gao and Wu ()] Qingtian Han; Xiaoyan Gao , ; Wenguo Wu . *9th International Conference on*
241 *Computer-Aided Industrial Design and Conceptual Design*, 2008. 2008. p. . (Study on Web Mining Algorithm
242 based on Usage Mining)
- 243 [Xiangjun and Shi Fufu ()] ‘Research of WEB Usage Mining Based on Negative Association Rules’. Dong
244 Xiangjun , ; Shi Fufu . *International Forum on Computer Science-Technology and Applications (IFCSTA*
245 *'09)*, 2009. 1 p. .
- 246 [Salin and Senkul ()] ‘Using semantic information for web usage mining based recommendation’. S Salin , P
247 Senkul . *24th International Symposium on Computer and Information Sciences*, 2009. p. . (ISCIS 2009)
- 248 [Zhang et al. ()] ‘Web Usage Mining Based On Fuzzy Clustering in Identifying Target Group’. Jianxi Zhang
249 , Peiyong Zhao , Lin Shang , Lunsheng Wang . *International Colloquium on Computing, Communication,*
250 *Control, and Management* 2009. 4 p. .
- 251 [Wu et al. ()] ‘Web Usage Mining on the Sequences of Clicking Patterns in a Grid Computing Environment’. Chih-
252 Hung Wu , Yen-Liang Wu , Yuan-Ming Chang , Ming-Hung Hung . *International Conference on Machine*
253 *Learning and Cybernetics (ICMLC)*, 2010. 6 p. .
- 254 [Hussain et al. ()] ‘Web usage mining: A survey on preprocessing of web log file’. T Hussain , S Asghar , N
255 Masood . *International Conference on Information and Emerging Technologies (ICIET)*, 2010. p. .
- 256 [Etminani et al. ()] ‘Web usage mining: Discovery of the users’ navigational patterns using SOM’. K Etminani ,
257 A R Delui , N R Yanehsari , M Rouhani . *First International Conference on Networked Digital Technologies*
258 *(NDT '09)*, 2009. p. .