

1 Web usage Mining: A Novel Approach for Web user Session 2 Construction

3 Neha Sharma¹ and Pawan Makhija²

4 ¹ SGSITS

5 *Received: 11 February 2015 Accepted: 28 February 2015 Published: 15 March 2015*

6 **Abstract**

7 The growth of World Wide Web is incredible as it can be seen in present days. Web usage
8 mining plays an important role in the personalization of Web services, adaptation of Web
9 sites, and the improvement of Web server performance. It applies data mining techniques to
10 discover Web access patterns from Web log data. In order to discover access patterns, Web log
11 data should be reconstructed into sessions. This paper provides a novel approach for session
12 identification.

14 *Index terms*— web mining, web server logs, web usage mining (wum), preprocessing, session identification.

15 **1 Introduction**

16 Web Usage Mining deals with understanding of user behavior, while interacting with web site, by using various
17 log files to extract knowledge from them. This extracted knowledge can be applied for efficient reorganization
18 of web site, better personalization and recommendation, improvement in links and navigation, attracting more
19 advertisement. As a result more users attract towards web site hence will be able to generate more revenue out
20 of it. [1]. Web Usage mining is made up with three procedures, as data preprocessing, data mining and pattern
21 analyzing. Data preprocessing contains three steps as data cleaning, user identification, session identification.
22 Session identification is an important step in data processing of web log mining. A session is defined as a group
23 of requests made by a single user for a single navigation. A user may have a single or multiple sessions during
24 a period of time. Presently sessions are identified either on Time based method or Navigation based method.
25 Here, we proposed a novel approach for user session identification by combining both Time based method and
26 Navigation based method.

28 **2 II.**

29 **3 Motivation**

30 Web log mining is to discover the mode of users' accessing to web page through mining web logs. In the process,
31 the designer's knowledge fields, the rate of his interesting and the users' visiting habit can be refined, which
32 can optimize the site's structure, develop individual service and the control of the users that is useful strategies
33 information for the designers and the managers. The most important and time-consuming link in mining web logs
34 is the session identification in web log preprocessing. The users' session is a session aggregation covering more
35 than one web services. The aim of session identification is to divide the users' page into an isolated identification.

36 **4 III.**

37 **5 Related Works**

38 The focus of literature review is to study, compare and contrast the available session identification techniques.
39 Traditional session identification algorithm is based on a uniform and fixed timeout. The set of pages visited by
40 a specific user at a specific time is called page viewing time. It varies from 25.5 minutes to 24 hours while 30

5 RELATED WORKS

41 minutes is the default timeout by Cooley [2]. If the interval between two sequential requests exceeds the timeout,
42 new session is determined.

43 Timeout algorithm uses a pre-fixed value of threshold for session identification in which if the interval between
44 two sequential requests exceeds the threshold value, a new session is determined. According to He Xinhua and
45 Wang Qiong [3], However, because of the uniform and fixed value, the algorithm cannot obtain efficient effect of
46 session identification in several situations like (1) Different user results different reading speeds, (2) Even by the
47 same user, different interest is shown on pages at different time, (3) Different page contains different contents.
48 Therefore, the time taken is often different. They propose a session identification algorithm based on dynamic
49 timeout, on the basis of traditional session identification algorithm. First, at beginning of the new session, the
50 initial timeout is set for each page using the formula? $0 = ? \cdot t \cdot (1 + ?)$

51 Where ? denotes smooth coefficient ranging from 1.1 ~1.6 and ? is an influence factor depend on link in and
52 link out of the page.

53 Second, while requested page is put into the current session, the timeout will be recomputed selectively in
54 order to make the timeout reflect the character of session using the formula? $= ? 0(t_{new} + t_0) / 2t_0$

55 Where t_0 denotes primal timeout of the page, and t_{new} denotes the timeout of the page that put intoW ©
56 2015 Global Journals Inc. (US)

57 Abstract-The growth of World Wide Web is incredible as it can be seen in present days. Web usage mining
58 plays an important role in the personalization of Web services, adaptation of Web sites, and the improvement of
59 Web server performance. It applies data mining techniques to discover Web access patterns from Web log data.
60 In order to discover access patterns, Web log data should be reconstructed into sessions. This paper provides a
61 novel approach for session identification.

62 current session and t_0 denotes the timeout by the adjustment last time.

63 In [7] Jozef Kapusta, Michal Munk and Martin Drlík, assume that the user goes over several navigation
64 pages during her/his visit until she/he finds the content page with required information. The content page is
65 a page where the user spends considerably more time in comparison with navigation pages. The content page
66 is considered, the end of the session. The division of pages into content and navigation pages is based on the
67 calculation of cut-off time C. When the cut-off time C is known the session can be created in such manner that
68 we compare the time of particular web page visit with the cut-off time C. The session is then defined as a path
69 through the navigation types of pages to the content page (the user spent there more time then C), they claim
70 the content page is last page of session. The cut-off time C is calculated on the basis of exponential distribution
71 of variable RLength(Time spent by user on individual page), here the assumption is that the variance of the
72 times spent on the auxiliary pages is small then the content page.

73 In [8] Zhixiang Chen, Richard H. Fowler and Ada Wai-Chee Fu, designed two algorithms for finding maximal
74 forward references (longest sequences of Web pages visited by a user without revisiting some previously visited
75 page in the sequence) from very large Web logs. They consider two types of sessions as ?interval session and ?
76 -gap session, where ?-interval session insures the duration of a session may not exceed a threshold of (30 minute)
77 and ?-gap session insures the time between any two consecutively assessed pages may not exceed a threshold of
78 ? (20 minute). They define a URL node structure to store the URL and the access time of a user access record
79 and a pointer to point to the next URL node and then the maximal forward reference session is calculated using
80 both interval session and gap session.

81 In [9] G. Arumugam, S. Sugana, Suggested algorithm which does not require searching whole tree representing
82 server pages. They employs concept of efficient use of data structure. Array List to represent web logs and user
83 access list, hash table for storing server pages, two way hashed structures for Access History List, represents user
84 accessed page sequences. Experiments reveals less time complexity and good accuracy of sessions generated as
85 compare to results of maximal forward reference method and reference length method.

86 In [10] Dr.Antony and V.Chitraa, proposed a new technique for identifying sessions for extraction of user
87 patterns. In the proposed method a matrix is constructed in which columns are the web pages and rows are
88 users. Browsing time (BT) for a particular page is determined by finding the differences between the time fields
89 of two consecutive entries of a same user and assumption is that the website Administrators fix the minimum
90 time and maximum time (BTmin and BTmax) for all web pages as per the contents. Codification of pages are
91 performed on the basis of BT, BTmin and BTmax and the sessions are calculated on the basis of this code. The
92 result is shown in the form of matrix.

93 In [11] Peng Zhu, Ming-sheng Zhao proposed an improved algorithm based on average time threshold value.
94 Experiments are conducted on the log files of Nanjing University Extra net user access logs. Because data of
95 log files are very large, They selected the log test algorithm of only one day ??March 15, 2008). Algorithm
96 proposed in this paper, takes individual differences into account to define the threshold value of users' browsing
97 pages, and identify long session page views, and divide the session less than the threshold into the next session.
98 They proposed two algorithms from which first algorithm constructs session of individual user and the second
99 algorithm disconnect the previous session into parts if there is no hyper link between two consecutive entries of
100 logs.

101 IV.

102 **6 Proposed Approach**

103 As we seen in the literature review the sessions are identifying either on the basis of time spend by user on
104 particular web page or on the basis of user navigation in web site topology.

105 Time based method ignores the web site structure, the sessions generated by such type of methods are not
106 generated right sessions as users reading speeds reflects the sessions. While in navigation method if particular
107 user not moves back, it not generates the sessions.

108 In our approach we combine both method to generate more informative sessions. Initially sessions are generated
109 by Maximal Forward reference method on these sessions the time based method have been applied with the
110 threshold value of 10 minutes. The experiment is conducted on the log data of www.smartsync.com of dated 8
111 Dec 2013.

112 V.

113 **7 Testing and Results**

114 The input data in this case are the access log files of the www.smartsync.com web server. Because data of log
115 files are very large, we select the log test dataset of only one day (dated 8 Dec 2013) of size 1 GB, 2 GB and 4
116 GB. Since the log data is very large in size it is not possible to count true sessions of whole data so we took
117 100 KB of data. In that data we manually found 53 true sessions and the number of session generated by the
118 existing methods and the proposed method are compared. For finding the accuracy of proposed approach we have
119 calculated the ratio of generated session and true session. Table-2 showing the comparison of existing method
120 and proposed method with true sessions. In the Table-2 S is number of sessions generated by methods and T is
121 the true sessions counted manually.

122 **8 Conclusion**

123 The growth of the web has resulted in a huge amount of information that is now freely offered for user access.
124 The several kinds of data have to be handled and organized in a manner that they can be accessed by several
125 users effectively and efficiently. The experiment on 4 GB data shows that the new method proposed in this report
126 generates more sessions (3102) than the traditional Time Based Method (2875) and Maximal Forward Sequence
127 Method (2742). On comparing with the true sessions on 100 KB data, the accuracy of session is increased to
88%.

Web usage Mining: A Novel Approach for Web user Session Construction
based
Method

Year	2. Maximal 968 Forward reference method	16402742
16		
Volume XV Issue III Ver- sion I	3. Proposed 1120 Approach	17103102
() E		
Global Journal of C omp uter S cience and T echnology		

© 2015 Global Journals Inc. (US) 1

Figure 1:

2

Methods	S	T	S/T %
Time based Method	42	53	79 %
Maximal Forward reference method	32	53	60 %
Proposed Approach	47	53	88 %
VI.			

Figure 2: Table 2 :

129 [Dr et al. (2011)] ‘A Novel Technique for Sessions Identification in Web Usage Mining Preprocessing’ Antony
130 Dr , V Selvadoss Thanamani , Chitraa . *International Journal of Computer Applications* November 2011. 34
131 (9) .

132 [Yuankang and Zhiqi ()] ‘A session identification algorithm based on frame page and page threshold’ Fang
133 Yuankang , Huang Zhiqi . *Computer Science and Information Technology (ICCSIT), 3rd IEEE International*
134 *Conference*, 2010.

135 [Huidrom and Bagoria (2013)] ‘Clustering Techniques for the Identification of Web User Session’ Nirmala
136 Huidrom , Neha Bagoria . *International Journal of Scientific and Research Publications* January 2013. 3
137 (1) .

138 [Kapusta et al. ()] ‘Cut-off Time Calculation for User Session Identification by Reference Length’ Jozef Kapusta
139 , Michal Munk , Martin Drlík . *IEEE* 2012.

140 [Xinhua and Qiong] *Dynamic Timeout-Based A Session Identification Algorithm*, He Xinhua , Wang Qiong .
141 *IEEE* 2011.

142 [Chen et al. ()] ‘Linear Time Algorithms for Finding Maximal Forward References’ Zhixiang Chen , Richard H
143 Fowler , Ada Wai-Chee Fu . *Intl Conf On Info Tech: Coding and Computing (ITCC03), Proc. of the*, 2003.
144 *IEEE*.

145 [Arumugam and Sugana ()] *Optimum algorithm for generation of user session sequences using server side web*
146 *user logs*, G Arumugam , S Sugana . 2009. *IEEE*.

147 [Zhu and Zhao] *Session Identification Algorithm for Web Log Mining*, Peng Zhu , Ming-Sheng Zhao . *IEEE* 2010.

148 [Zhang and Ghorbani ()] ‘The Reconstruction of user session from a server log using improved time oriented
149 heuristic’ J Zhang , Ali A Ghorbani . *IIInd Annual Confernnce on Communication Networks and Service*
150 *Research* 2004. *IEEE*.

151 [Robert et al. ()] ‘Web mining: Information and Pattern Discovery on the World Wide Web’ Robert , Bamshed
152 Cooley , Jaideep Mobasher , Srinivastava . *International conference on Tools with Artificial Intelligence*,
153 (Newport Beach) 1997. *IEEE*. p. .

154 [Chintan et al.] *Web Usage Mining: A Review on Process*, R Chintan , Nirali N Varnagar , Madhak , M Trupti
155 , Jayesh N Kodinariya , Rathod . *IEEE* 2013.

156 [Dell ()] ‘Web user session reconstruction using integer programming’ R F Dell . *International Conference on*
157 *Web Intelligence and Intelligent Agent Technology*, 2008.