Artificial Intelligence formulated this projection for compatibility purposes from the original article published at Global Journals. However, this technology is currently in beta. *Therefore, kindly ignore odd layouts, missed formulae, text, tables, or figures.* 

1 2	Web Page Recommendation Approach Using Weighted Sequential Patterns and Markov Model
3	Dr. K. Suneetha <sup>1</sup> and Dr.M.Usha Rani <sup>2</sup>
5	Received: 13 December 2011 Accepted: 3 January 2012 Published: 15 January 2012

#### 7 Abstract

20

Web page recommendation aims to predict the user's navigation through the help of web 8 usage mining techniques. Currently, researchers focus their attention to develop a web page 9 recommendation algorithm using the well known pattern mining techniques. Here, we have 10 presented a web page recommendation algorithm using weighted sequential patterns and 11 markov model. To mine the weighted sequential pattern, we have modified the prefixspan 12 algorithm incorporating the weightage constraints such as, spending time and recent visiting. 13 Then, the weighted sequential patterns are utilized to construct the recommendation model 14 using the Patricia trie-based tree structure. Finally, the recommendation of the current users is 15 done with the help of markov model that is the probability theory enabling the reasoning and 16 computation as intractable. For experimentation, the synthetic dataset is utilized to analyze 17 the performance of W-Prefixspan algorithm as well as web page recommendation algorithm. 18 From the results, the memory required for the W-prefixSpan algorithm is less than 50 19

Index terms— Web page recommendation, Weighted sequential pattern, Prefixspan, Patricia-trie, Markov model.

## 23 1 Introduction

equential pattern mining, an advance of association rule mining, is an imperative subject of data mining, often 24 applied for extracting the useful information [9]. Sequential pattern mining algorithms deals with the problem of 25 determining the frequent sequences in a given database [6]. Sequential pattern is a sequence of itemsets that often 26 27 occur in a specific order and thus, all items in the same itemset are expected to encompass the same transaction time value or within a time gap. Each sequence have a temporally ordered list of events, wherein each event is 28 a compilation of items (itemset) occurring simultaneously. The temporal ordering among the events is induced 29 by the absolute timestamps associated with the events [10]. Generally, all the transactions of a customer are 30 collectively viewed as a sequence, called customer -sequence, where each transaction is represented as an itemset 31 in that sequence and all the transactions are scheduled in a particular order based on the transaction-time [8]. 32 Recently, there has been a substantial interest in using sequential mining approaches to construct web page 33 34 recommendation systems.

Here, we have presented a web page recommendation algorithm using weighted sequential patterns and markov model. The overall process of page recommendation based on Web usage mining consists of three phases: data preparation, mining of weighted sequential patterns and recommendation.

Data preparation: This step consists of, (i) identifying the different users' sessions from the usually very poor information available in web log files and (ii) reconstructing the users' navigation path within the identified sessions.

41 Mining of weighted sequential patterns: The traditional sequential pattern mining problem is extended by 42 allowing a weight to be associated with each page in a user session to reflect interest of each page within the user 43 session. In turn, this provides with an opportunity to associate a weight parameter with each page in a resulting
44 sequential pattern, which called a weighted sequential pattern (WSP).

45 Recommendation: After mining the weighted sequential patterns, the patricia-based tree is constructed. From 46 the Patricia tree, a recommendation model is developed based on markov model for predictions of users to find 47 web pages they want to visit.

### 48 **2** II.

## <sup>49</sup> **3** Related work

Literature presents a lot of web page recommendation algorithms based on web used mining, collaborative filtering, and rule-based filtering. Here, we portray some recent works presented in the literature for web page recommendation.

Forsati, R et al [3] have developed an effective and scalable approach to deal with the web page recommendation 53 problem. Here, a distributed learning machine has been employed to learn the performance of previous users' 54 55 and to cluster the pages based on learned pattern. Dealing with unvisited or recently added pages was one 56 of the difficult and challenging tasks in recommendation systems. As they would never be recommended, it is indispensable to provide a chance for these seldom visited or recently added pages to be incorporated in 57 58 the recommendation set. By considering this problem, a Weighted Association Rule mining algorithm has 59 been presented for the recommendation purposes. Also, a HITS algorithm has been exploited to extend the recommendation set. Furthermore, they have analyzed the proposed algorithm under various settings, and 60 revealed the efficiency of this approach in enhancing the overall quality of web recommendations. 61

Yicen Liu et al [15] have introduced an automatic tag recommendation algorithm that has been employed 62 in the large-scale and real-time data process successfully and efficiently. Most of the prior researches on tag 63 suggestion have focused on initially, discovering the relationship between testing and training data and then, 64 65 assigning the top ranked tags of the most related training data to the testing object. But, they not paid any 66 attention in determining the internal relationship between the tags and weblogs. In their current research, more than 43% of tags, which have been labeled by weblog users, have really been employed in the body of the text. 67 68 In the meantime, the term frequency distribution, the paragraph frequency distribution, and the first occurrence position of tags were dissimilar from the ones of non-tags in the text. As well, the tags of a weblog have been 69 assigned in two steps. Initially, some probability distributions of the word attributes have been trained through 70 the labeled training weblogs, and some keywords of a testing weblog have been extracted as one part of the tags 71 based on the probability distributions. Subsequently, with the aid of Latent Semantic Indexing (LSI) model, 72 the other parts of the tags have been obtained from the first part ones. Experiments conducted on an extensive 73 74 tagging dataset of weblogs 12 have confirmed that the average tagging time for a new weblog was less than 0.0275 sec, and more than 74% of testing weblogs have been properly labeled by means of the top 15 tags.

76 An approach for recommendations of unvisited pages has been presented by Forsati, R et al [11]. They have focused on the recommender systems based on the user's navigational patterns and provided proper 77 78 recommendations to cater to the current needs of the user. The group of users with analogous browsing patterns has been identified by employing an offline data preprocessing and clustering technique. The experiments 79 conducted on real usage data from a commercial web site have demonstrated a considerable enhancement in 80 the recommendation efficiency of the proposed system. Web Personalization is viewed as an application of data 81 mining and machine learning approaches to create models of user behavior that can be applied to the task of 82 forecasting the user needs and adapting future interactions with the eventual goal of enhanced user satisfaction. 83 84 An extensive overview of intelligent methods for Web Personalization has been presented by Sarabjot Singh 85 Anand and Bamshad Mobasher [12]. They have studied the state-of-the-art in Web personalization.

Initially, a depiction of the personalization process and a classification of the current techniques to Web personalization have been presented. Also, they have discussed the different sources of data available to personalization systems, the modeling techniques utilized, and the current techniques to analyze these systems. Numerous challenges faced by the researchers in developing these systems and also the solutions to these challenges proposed in literature have been described. They have concluded with a discussion on the open challenges that must be addressed by the research community if this technology is to create a positive impact on user satisfaction

92 with the Web.

Due to the increasing number of Web sites such as e-businesses contain a huge number of pages, users find it 93 very hard to swiftly reach their own target pages. Thus, Hiroshi Ishikawa et al [5] have proposed two approaches 94 95 to Web usage mining as a key solution to these problems. First of all, an efficient recommendation system 96 called the L-R system has been depicted, which creates user models through classifying the Web access logs 97 and by mining access patterns based on the transition probability of page accesses, and then, recommend the 98 significant pages to the users based on both the user models and the Web contents. The prototype system has been analyzed and obtained the positive effects. Secondly, another approach has been employed for creating user 99 models that clusters the Web access logs based on the access patterns. Moreover, the user models assist to find 100 the unexpected access paths corresponding to ill-formed Web site design. In addition, Daniel Mican and Nicolae 101 Tomai [2] have proposed WRS, architecture for robust web applications. Within the structure, usage data was 102 being implicitly obtained by data collection sub-module. Here, the usage data has been extracted, online and in 103

real time, via a proactive technique. They have efficiently exploited association rule mining among both frequent
and infrequent items for successful pattern discovery. This was due to the fact that the pattern discovery module
transactionally processes users' sessions and employs incremental storage of rules. Also, they have proved that
the Wise Recommender System (WRS) has been straightforwardly implemented within any web application,
because of the efficient integration of the three phases into an online transactional process.

# <sup>109</sup> 4 Global Journal of Computer Science and Technology Volume <sup>110</sup> XII Issue IX Version I

111 III.

# <sup>112</sup> 5 Proposed algorithm to web page recommendation based on <sup>113</sup> weighted sequential access patterns

Web page recommendation is significant research over the past decade due to its real world application. With the intention of real world applicability, we have developed an approach for web page recommendation using weighted sequential pattern and markov model. Here, the traditional sequential pattern mining algorithm is modified significantly by incorporating the significant measure to mine more useful patterns. Then, the markov model described in [4] is used to recommend the web pages. The block diagram of the proposed approach for web page recommendation is given in Figure 1. The important steps for generating recommendations to the user is as follows,

? Data Preprocessing This section describes the preprocessing steps of web log file that is the input of 121 the proposed web page recommendation approach. In general, the web log file consists of, IP address, access 122 time, HTTP request method used, URL of the referring page and browser name (for an example, Web server 123 log file: 192.162.37.21 [23/Feb/2012:08:17:25]"GET / HTTP/1.1" "http://www.sigkdd.org/kddcup/index.php" 124 IE/7.0 Windows 07). The initial process of the proposed web page recommendation approach is to preprocess 125 the web log file such a way that the mining process should be applied. Here, we make use of the sequential 126 pattern mining process so there is a need to convert the web log file into the sequential database that should be 127 in the proper format to mine the weighted sequential patterns. 128

User Identification: User identification is an important step for constructing the sequential database. IP address and user session are utilized here to track a unique user from the web log file. Unique IP address is a new user but at the same time, the user session should be fixed for particular time period. If the user session is reached to a particular duration for the same IP address, then the new session is acted like new user. Based on this, user transaction is formed from the web log file.

134 Weighted Sequential database generation:

Once the user transactions are identified, the weighted sequential database is generated including the sequence of web pages visited by the user, time spent by the user on corresponding web page and its recent information.

Let assume the web log database D having IP address, access time, HTTP request method used, URL of the referring page and browser name. After applying the data preprocessing steps, we generate the weighted sequential database that is represented as follows, ij W in which, ' i ' belongs to the set of users and ' j ' signifies the set of pages visited by the corresponding user. Here, every element of ij w contains three tuples ), , (ij ij ij iii ij r s p w =

in which the first tuple belongs to the web pages, the second one belongs to the time spent within that page and the third one belong to whether its recent one or not. For example, 0, 0, 20, (1 p w ij = is a tuple, whichdenotes that the user spent 20 seconds on page 1 p and '0' signifies that the page 1 p is not recently accessed by the user.

# <sup>146</sup> 6 b) W-Prefixspan for Mining of Weighted Sequential Web <sup>147</sup> Access Pattern

Once the weighted sequential database is constructed, the mining procedure is carried out to find the interesting 148 sequential patterns. Here, PrefixSpan [7] is modified as W-PrefixSpan (Weighted-PrefixSpan) by incorporating 149 the spending time and recent view into the mining procedure. The weightage measure assumed in the proposed 150 W-PrefixSpan algorithm is spending time and recent view. The two aspects taken for providing the weightage 151 152 of the sequential patterns are, Spending time: One of the fields in the web log data is the duration of the web page which is viewed by the user. Generally, time spent by the user within a particular page is necessary to 153 154 identify the importance of web pages. From the web page which having long duration, we can conclude that this 155 particular web page has been referred by the user in a long occasion because of its worth. Thus, the spending time is an important measure for the researchers who are attempting to identify the interest of the users. So, if we 156 incorporate the time duration into the mining of sequential patterns, the interesting relationships can be found 157 out from the mined sequential patterns that can be effectively applied to web page recommendation process. 158 Recent view: Another significant measure taken for sequential pattern mining is recent view that describes 159

whether the page is accessed recently or not. The reason behind taking the recent view for mining the sequential

#### 10 A) EXPERIMENTAL SET UP AND DATASET DESCRIPTION

pattern is that the more importance should be given for the web pages which are accessed recently than the older one. The behavior of the user surely vary depend on the time so the recent behavior of the user is significant for finding the sequence analysis. With the intention of behavior variation over time, the recent view is also incorporated into the sequential pattern mining algorithm to achieve a subset of more interesting sequential patterns (SR-Patterns).

W-PrefixSpan algorithm: In this section, we describe an efficient algorithm, W-PrefixSpan, for mining all the SR-patterns from weighted sequence databases. The W-PrefixSpan algorithm is developed by modifying the eminent PrefixSpan algorithm, which uses the pattern growth methodology for mining the frequent sequential patterns recursively. Let ), ,(,), , , (), , , (1 1W only if, (1) s W is a subsequence of ij W, s ij W W? (2) m t t t < < < ? 2 1

where, 1 t is the time Where, T M ? Total number of web pages in one transaction, i s ? spending time Then, the projection database is formed by projecting the collection of postfixes of mined 1-SR sequence. In projection database, 'n ' disjoint subsets are generated if the mined 1-SR patterns contain 'n ' number of sequence.at which ij p occurred in s W , m r ? ? 1 . A sequence is said to be SR sequence ij W if and only if, (1) s W is a subsequence of ij W , (2

Then, the 2-length SR-patterns are mined from the projected database by computing the weighted support on the projected database. Again, the projected database is formed with the help of mined 2-SR patterns and this process is repeated recursively until all SR sequential patterns are mined. The following provides a detailed explanation of the important steps involved in the proposed W-PrefixSpan algorithm. Method: 1. Scan? | ij W once, find the set of SR-pattern s such that a) 's ' can be assembled to the last element of ? to form a SR-sequential pattern; or b) s can be appended to ? to form a SRsequential pattern.

2. For each SR-web page s, append it to ? to form a SR-sequential pattern '?.

5. For each '? , construct '? -projected database ? | ij W , and call PrefixSpan ( , 1 , ' + l ? ? | ij W ).

# <sup>185</sup> 7 Global Journal of Computer Science and Technology Volume <sup>186</sup> XII Issue IX Version I c) Building of Pattern Tree Model

Once we mine the weighted sequential patterns, the pattern tree is constructed using the procedure defined 187 in [14,13]. Initially, trie-based data structure given in [1] is used to construct the pattern tree for web page 188 recommendation. Later, the modification was done by [14,13], who utilized the patricia-based data structure 189 for web page recommendation due to the advantages of particia structure over the trie structure. Here, the 190 procedure defined in [14,13] is applied to the proposed approach for constructing tree structure. The method for 191 constructing the pattern tree in the proposed web page recommendation approach is as follows: 1) Generate an 192 empty root node, 2) Add the most sub pattern in the SR-sequential pattern set into a node next to the root node, 193 3) Insert the postfixes of pattern into child node only if the current pattern to be inserted is a super pattern of 194 inserted patterns, 4) Otherwise, current pattern is inserted into the node next to the root node, and 5) Step 3 195 and step 4 is repeated for every pattern in the mined SR-pattern set. 196

### <sup>197</sup> 8 d) Generation of Recommendations Using Markov Model

This section describes the markov model utilized for web page recommendation. Here, we make use of the markov model described in [4] that is used in the identication of the next page to be accessed by the Web site user based on the sequence of previously accessed pages. Here, whenever a new user comes to get the recommendation, the sequence path of the new user is matched with the patricia-trie structure. Then, the subsequent web page whether it may be from same node or from its child node is retrieved. Now, the sequence path of the new user is used to find the accurate recommendation using the probability definition used in the previous work [4]. Let the input sequence visited by the user be

### 205 9 Results and discussions

This section presents the results obtained from the experimentation and its detailed discussion about the results. The proposed approach of web page recommendation is experimented with the synthetic dataset and the result is evaluated with the precision, applicability and hit ratio.

is evaluated with the precision, applicability and hit ratio.

## <sup>209</sup> 10 a) Experimental Set Up and Dataset Description

The proposed web page recommendation approach is implemented in Java (jdk 1.6). Here, the synthetic dataset is generated as like the same format of real datasets and the performance of the proposed approach is evaluated with the evaluation metrics. The generated synthetic dataset is divided into two parts such as, Training dataset (It is used for building the pattern tree model and test dataset (It is used for testing the web recommendation approach). The figure **??** shows that the execution time of the W-prefixSpan algorithm is 1.5 Sec that is less compared with the time taken of the PrefixSpan algorithm. In memory usage, the W-prefixSpan algorithm needed only less than 50% memory compared with Prefixspan algorithm. Figure **??** shows that number of patterns mined

217 using W-Prefixspan is high compared to Prefixspan algorithm.

218 V.

## 219 11 Conclusion

We have proposed a web page recommendation algorithm using weighted sequential patterns and markov model. Here, we have presented W-PrefixSpan algorithm that is developed by incorporating the weightage constraints such as, spending time and recent visiting with the prefixspan algorithm. The mined weighted sequential patterns are then utilized to construct the recommendation model using the patricia trie-based tree structure. At last, markov model-based recommendation is carried out for the current users by matching the visiting path with the

225 tree and markov model. The experimentation is done with the help of synthetic dataset and the performance

of W-Prefixspan algorithm as well as web page recommendation algorithm is analyzed. From the results, the memory required for the W-prefixSpan algorithm is less than 50% of memory needed for PrefixSpan algorithm.



Figure 1: Fig. 1:



Figure 2:



Figure 3:



Figure 4:

227 228 1 2 3 4 5

 $<sup>^1 \</sup>odot$  2012 Global Journals Inc. (US) 2012 April

 $<sup>^2 @</sup>$  2012 Global Journals Inc. (US) Global Journal of Computer Science and Technology Volume XII Issue IX Version I

 $<sup>^3 @</sup>$  2012 Global Journals Inc. (US) 2012 April

 $<sup>^4 \</sup>odot$  2012 Global Journals Inc. (US) 2012 April

 $<sup>^5 @</sup>$  2012 Global Journals Inc. (US)

- [Orlando et al. ()] 'A new algorithm for gap constrained sequence mining'. Salvatore Orlando , Raffaele Perego
   , Claudio Silvestri . Proceedings of the ACM Symposium on Applied Computing, (the ACM Symposium on
- 231 Applied ComputingNicosia, Cyprus) 2004. p. .
- [Hou and Zhang ()] 'Alarms Association Rules Based on Sequential Pattern Mining Algorithm'. Sizu Hou ,
   Xianfei Zhang . proceedings of the Fifth International Conference on Fuzzy Systems and Knowledge Discovery,
   (the Fifth International Conference on Fuzzy Systems and Knowledge DiscoveryShandong) 2008. 2 p. .
- [Forsati et al. ()] 'An efficient algorithm for web recommendation systems'. R Forsati , M R Meybodi , A Rahbar
   *IEEE/ACS International Conference on Computer Systems and Applications*, 2009. p. .
- [Utpala Niranjan et al. (2010)] 'An Efficient System Based On Closed Sequential Patterns for Web Recommen dations'. R B V Utpala Niranjan , V Subramanyam , Khana . *IJCSI International Journal of Computer Science Issues* May 2010. 7 (4) .
- [Zhou et al. ()] 'An Intelligent Recommender System using Sequential Web Access Patterns'. Baoyao Zhou , Siu
   Cheung Hui , Kuiyu Chang . proceedings of IEEE Conference on Cybernetics and Intelligent Systems, (IEEE
   Conference on Cybernetics and Intelligent Systems) 2004.
- [Mican and Tomai ()] 'Association-Rules-Based Recommender System for Personalization in Adaptive Web Based Applications'. Daniel Mican , Nicolae Tomai . Proceeding ICWE'10 Proceedings of the 10th international
   conference on Current trends in web engineering, (eeding ICWE'10 eedings of the 10th international conference
- on Current trends in web engineering) 2003. p. .
- [Sumathi et al. ()] 'Automatic Recommendation of Web Pages in Web Usage Mining'. C P Sumathi , R Valli ,
  T Santhanam . IJCSE) International Journal on Computer Science and Engineering 2010. 02 (09) p. .
- [Liu et al. ()] 'Automatic Tag Recommendation for Weblogs'. Yicen Liu , Mingrong Liu , Xing Chen , Liang
   Xiang , Qing Yang . International Conference on Information Technology and Computer Science, 2009. p. .
- [Utpala Niranjan et al. ()] 'Developing a Web Recommendation System Based on Closed Sequential Patterns'. R
   B V Utpala Niranjan , V Subramanyam , Khanaa . Communications in Computer and Information Science
   2010. 101 (1) p. .
- [Khalil et al. ()] 'Integrating Recommendation Models for Improved Web Page Prediction Accuracy'. Faten Khalil
   Jiuyong Li , Hua Wang . Proceedings of the thirtyfirst Australasian conference on Computer science, (the
   thirtyfirst Australasian conference on Computer science) 2008. 74.
- [Singh Anand and Mobasher ()] 'Intelligent Techniques for Web Personalization'. Sarabjot Singh Anand ,
   Bamshad Mobasher . *ITWP 2003, LNAI 3169*, (Berlin Heidelberg) 2005. Springer-Verlag. p. .
- [Jatin et al. (2007)] 'Modified Web Access Pattern (mWAP) Approach for Sequential Pattern Mining'. D Jatin ,
   Sanjay Parmar , Garg . Journal of computer Science June 2007. 6 (2) p. .
- 261 [Pei et al. ()] 'PrefixSpan,: mining sequential patterns efficiently by prefix-projected pattern growth'. Jian Pei ,
- Jiawei Han, B Mortazavi-Asl, H Pinto, Qiming Chen, U Dayal, Mei-Chun, Hsu. Proceedings of 17th International Conference on Data Engineering, (17th International Conference on Data Engineering) 2001.
- 264 [Ishikawa et al. ()] Proceeding Revised Papers from the NODe 2002 Web and Database-Related Workshops
- on Web, Web-Services, and Database Systems, Hiroshi Ishikawa , Manabu Ohta , Shohei Yokoyama ,
- Junya Nakayama , Kaoru Katayama . 2003. p. . (On the Effectiveness of Web Usage Mining for Page Recommendation and Restructuring)
- [Zhao and Bhowmick ()] 'Sequential Pattern Mining: A Survey'. Qiankun Zhao , Sourav S Bhowmick . CAIS
   2003. (118) . Nanyang Technological University (Technical Report)