# Web Usage Mining for Comparing User Access Behaviour using Sequential Pattern

Amit Dipchandji Kasliwal[#], Dr. Girish S. Katkar[*]

[#]*Malegaon, Nashik, Maharashtra, India*
[*]*Dept. of Computer Science, Arts, Sci., Commerce College,*
*Koradi, Nagpur, Maharashtra, India*

*Abstract -* **In the current era of technology, every organization using a dynamic process of communication that changing in very interactive manner which is fulfilled by their own website hosted by their own web servers or commercial web server. This website gathers information about user access at every time when user interacts for communication with some resources available to them. And this makes easiest way for website administrator to compare the information about the user's navigational patterns explored from the web access logs and can also be useful to compare the user access behaviour about what will be the next page that should be provided to the user following the patterns. To achieve this purpose sequential pattern technique can be used. Sequential pattern technique discovers frequent sequence referred to as patterns in sequential database. In this study we are proposing the algorithm for finding the sequential pattern that could suggest and compare the generated patterns with the minimum support measure.**

*Keywords—* **Web Usage Mining, Web access log file, Sequential patterns.**

## I. INTRODUCTION

Sequential pattern technique discovers frequent sequence referred to as patterns in sequential database. The sequential pattern mining process was introduced by Agrawal et al and Srikant et al [1] based on their study regarding the customer purchase sequences. The problem they studied under process was to find all frequent sequences such that sequences in which items occurs frequently in the set of sequences is not less than the minimum support.

An access made to each webpage of website, this access is get recorded into web access log file. Web log file can be considered as a sequence of web pages maintained in the form of log entry accessed by users. As the main intension of using web usage mining is to discover interesting and frequent user access patterns from the web access log data as shown by Pei et al [9]. Web access log stores the information that can be important for website administrator while improving web sites by link to link access by Srivastava et al [1]. In this chapter the authors are proposing and implementing SequentialPattern algorithm in WebUMining model for comparing user access behaviour. The goal of this chapter is to find sequential patterns from web log file that occurs frequently with respect to minimum support count provided by the user. As, it is an important data mining task and has broad applications, the proposed algorithm is implemented to provide the requested user with

more useful and interesting patterns. A web access by user is a pattern that is pursued frequently by users. Using these sequences as starts, researcher is converting it into the sequential database so that it can be mine with proposed algorithm. The proposed algorithm mines the complete set of patterns but greatly reduces the efforts of candidate generated through subsequence from patterns. Next on, the proposed work substantially reduces the size of irrelevant sequential patterns and hence leads to efficient processing.

Sequential pattern mining is an important data mining technique can be used for retrieving patterns from large size of data. It is also used to mine frequent navigational paths among user accesses. Web Access Pattern is a sequential pattern technique in a large collection web log data. A web log record can be referred as the sequence of pairs containing user's IP id and requested URL. Sequential pattern mining is an efficient technology used for extracting access patterns. For the sequential pattern process, let us consider itemset $I_s = ( i_1, i_2, …, i_n)$ and $X$ is the subset such that $X \subset I_s$. Let $S$ is the sequence in ordered manner list of itemset denoted by $(S_1 S_2 ... S_i)$, where $S_i$ is itemset of the sequence denoted by $(X_1 X_2 ... X_m)$ where $X_m$ is the item in sequence $S_i$. If an item occurs only one item then it can be written as $X$. The number of times the item occurs in a sequence is called the length of the sequence. A sequence is then get store into database $D$ as a set of element as $\{S_{id}, S\}$, where $S_{id}$ is sequence identifier and $S$ is sequence. The support count for a sequence $\alpha$ in sequential database $D$, $S$ is the number of tuples in the database containing $\alpha$, given as $SUPPORT_S (\alpha) = | \{S_{id}, S\} \ni (S, \alpha \in D ) |$. For the user defined support count, a sequence $\alpha$ is said to be called a sequential pattern in sequential database $S$ if and only if $SUPPORT_S(\alpha) \geq Min\_Support\_Count$. To demonstrate the web access pattern through web log access data, let us consider set of user access during a particular session denoted by $E$ used to represent web pages accessed by users. A web access sequence is ordered sequence that gives the ordered list web pages accessed by user in a session. It is given by S as $S = \{(E_1, E_2, ..., E_n) | E_i \ni E$ and $1 < i < n \}$ and in the sequential pattern the $E_i$ and $E_j$ are not necessary to be different.

Sequential Pattern Mining Algorithm is basically based on two approaches apriori based and pattern growth based approach. It mines frequent patterns and sequential pattern containing frequent itemset which are treated uniformly in

the training set. For example, in sequential pattern mining, a sequential pattern {(computer, software), (software, antivirus)} can be discovered with user defined support count as it is relatively high, almost in all scenario. And that makes it's limitation for the real life itemset. If we consider the real life scenario, items have important sequences such as {(event, charity), (car, insurance), (computer, webcam)} could not be mine with the traditional sequential pattern mining approach because the itemset contains item that may have low support count. And hence, in real life scenario items with less (or low) support count become more important due to some features of the item may have. After understanding this scenario, the proposed SequentialPattern algorithm is introduced by implementing the weight based pattern mining algorithm again with the modification in the main algorithm suggested by Chen et al, Huang et al [11]. In traditional sequential pattern mining technique, all sequences basically considered with same importance, but as in real life examples, sequences varies in their importance.

The proposed algorithm implemented here is with the key concept that to assign some weightage to the items according to their importance as they occurs in the sequence. In the consequence to the word weightage, the word weight can be defined as the value assigned to each item from itemset with respect to their importance. A pair $(X_i, W_i)$ is called a weighted item where $X_i \in X$ is an item and $W_i \in W$ is the weightage associated to $X_i$ where $W_i$ is a set of positive integer assigned to an item and $X=(X_1, X_2,..,X_n)$ be the set of unique items. A record is a set of weighted items in which a single item may appear in multiple record with different weightage associated with it. In this way the weightage is assigned to the items that belongs to the frequent item set prepared either by apriori algorithm or pattern growth approach.

## II. PROPOSED SEQUENTIALPATTERN ALGORITHM

In Sequential Pattern technique for mining, there are two module in which the first is to assign the weightage to the each frequent item and then discovering sequential pattern. In the proposed SequentialPattern algorithm, the researcher implementing these two parts where in first part, the occurrence of user accesses to webpage is used as to assign the weightage to each item from an access sequence. In the second part of the proposed algorithm, pattern discovering algorithm considers the weight of sequence and used it for finding out the support determined for user session and depending upon this the sequential patterns are discovered. The proposed SequentialPattern algorithm uses modified form of the structure used in WAPTree algorithm proposed by Pei et al[9] to deal with user access. In the tree generated with WAPTree algorithm, nodes of the tree are modified to handle the weight assigned to the items as the items are now the webpages. In the second part of the proposed SequentialPattern algorithm, the researcher uses the data structures and modified it as proposed in FOLMine algorithm by Rajimol et al [7]. These data structures used as per the data available in the form of web access log data for the proposed algorithm as listed below:

| Modified Data Structure | Description |
|---|---|
| **Item_List (ItLt)** | Linked list modified for containing the items and their occurrence in the web access log data. |
| **Curr_Item_List (CrItLt)** | Linked list modified so that it can be used for storing the items from the current user access webpages. |
| **First_Occ_List (FsOcLt)** | Linked list modified for storing the first occurrences of a given item in the web access log data. |
| **Acc_List (AcLt)** | Linked list modified for storing each web access sequence. |
| **Header_List[]** | Array structure containing start address of list containing item with their weight represented as Header_List. |

Table 1: Data Structures used in the proposed algorithm

## III. PSEUDO CODE AND FLOWCHART

The proposed algorithm at first scans each web access log entry from web access log one after and then stores it with assigning weightage of each item to 0. While reading the items from itemset, the Curr_Item_List get is loaded with items those are present in the access sequence and their occurrences. After completing the scanning of access sequence from web access log, the Item_List is get updated using Curr_Item_List. Perhaps, the weightage of each item from Curr_Item_List is calculated by performing division between the occurrence and total length of the access sequence and weightage in the Acc_List is get updated. After this preprocessing get over, the frequent itemset of weighted items is generated using Item_List. If the occurrence for an item is goes more than support, it is considered as the perfect sequential pattern for which website administrator is looking for.

As the data collection and data preprocessing task is already performed and discussed in chapter 2 hence researcher is starting the algorithm with step where WebUMining model is accepting the input file from user. This proposed algorithm is provided with the following pseudo code hence the researcher is intends to have the implementation of the proposed algorithm in the low level prototype model.

Algorithm Name: **SequentialPattern**
Input: Sequential data from Web Access Log,
Minimum Support,
Minimum Confidence.
Output: Sequential Pattern with Support and Confidence.
1. Begin
2. While eof (SequentialWebAccessLogFile)
   Read each web access log AND each item *i=0*
   Set length = 0
   For each item *i* in *S*
       Do length + +
           If (item *i* is not in Curr_Item_List)
               make entry in Curr_Item_List for this *i*
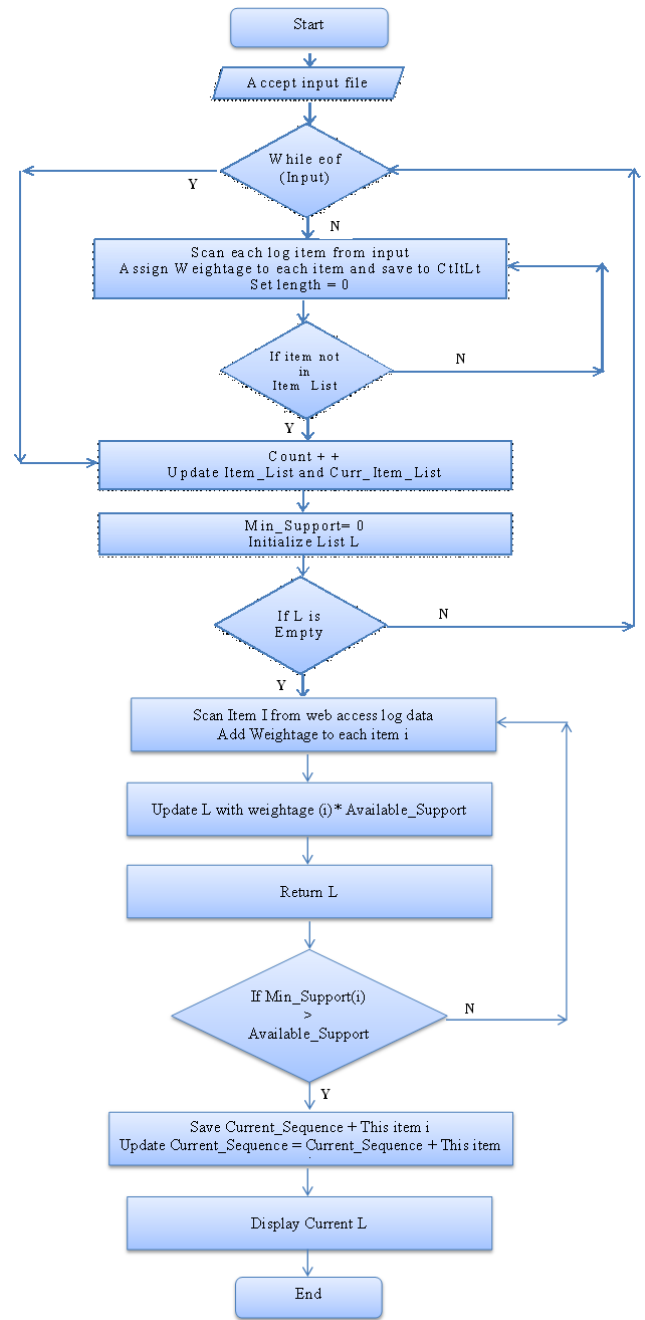
with   count($i$) = $1$

    Else

        count ($i$) ++

    End If

End For

Update the Item_List with the Curr_Item_List

For each $i$

    Update weight ($i$) = count ($i$ ) / length

End While

3.  For weighted item, $i$

    Set Min_Support = 0 and initialize $List_1$

4.  If $L_1$ is empty

    Then find first occurrence of $i$

    Else If

        *GOTO* step 2 to locate the first occurrences of the element $i$

    End If

5.  Add weightage to item for each occurrence

6.  Update $L_1$ with total weightage * Min_Support

7.  Return $L_1$

8.  If (Min_Support($i$) > Available_Support)

        Then Current Pattern + This Item $i$

        Save Current Pattern + This Item $i$ to Sequences

        Current Pattern = Current Pattern + This Item $i$

        *GOTO* step 5

    *End If*

9.  Display current $L_1$

10.  *End*

While implementing the proposed SequentialPattern algorithm, initially there was no method to follow a particular session for target page. And hence, the last page in session is taken as the target page. As each webpage refreshed, causes the duplication of webpages in the sequence, therefore nearby duplicate requests from single webpage are eliminated. The output of the proposed SequentialPattern algorithm is in the form of a text file containing the list of sequential patterns with their occurrences as well as support measure and confidence measure. In this format the results are then be easily represented by visualization software. Effectiveness of the proposed SequentialPattern algorithm is implemented and evaluated through comparing processing time and number of patterns generated during the web access log data submitted as as input to algorithm.

Following is the flowchart of the proposed SequentialPattern algorithm as discussed above. The flowchart contains data structures as discussed. The flowchart gives the idea about the flow of data and the actual process implemented in the proposed SequentialPattern algorithm. The algorithm is then implemented with JAVA and it is one of the algorithms of WebUMining model used for web usage mining for comparing ing user access behaviour.



## IV. EXPERIMENT

In this section, authors worked on the web access log data to perform the experiment to evaluate the sequential pattern that can be used for comparing user access behaviour through the proposed algorithm implemented in WebUMining model. During experiment, the focus is on fining sequential patterns from web log data, then looking for the presence of the weighted sequential pattern. For conducting the experiments researcher used the web log access data collected from kdnugget's repository named 'kddaccess.arff' which is preprocessed with the DataCleaningPreprocessing algorithm as discussed in chapter 4. Then this file is converted from text file to arff file format using WEKA as shown in figure 1 given below:

Figure 1: Input web access log data file

The WebUMining model accepts this input file and selecting the proposed SequentialPattern algorithm the remaining processes are shown in the following figures.
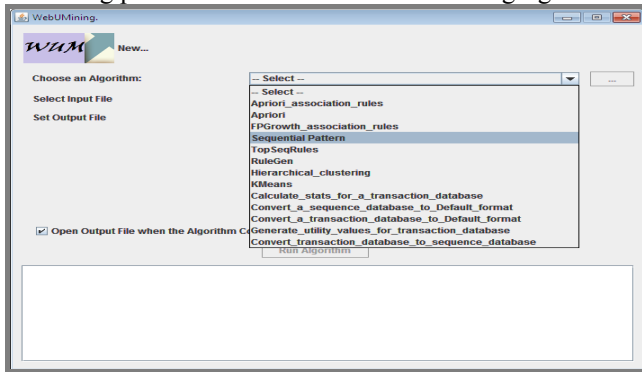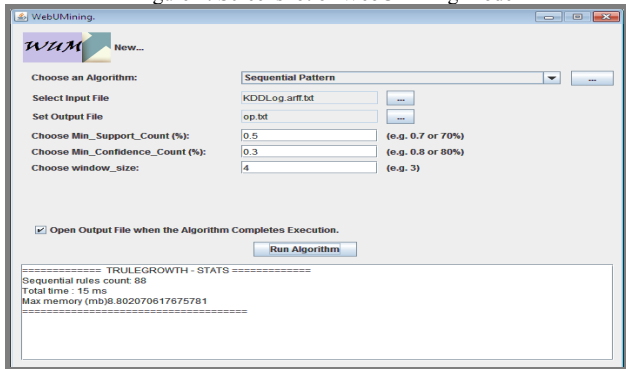

Figure 2: Screenshot of WebUMining Model
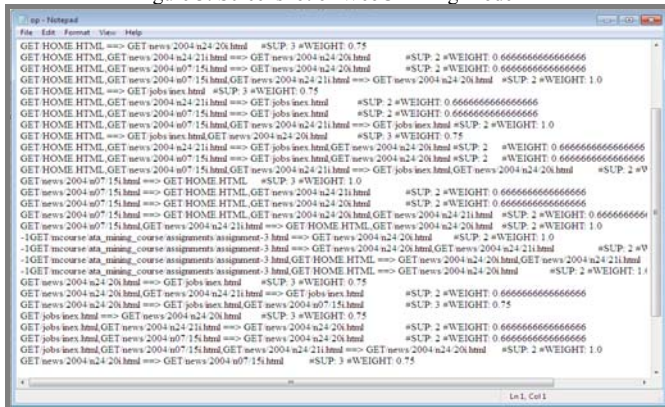

Figure 3: Screenshot of WebUMining Model

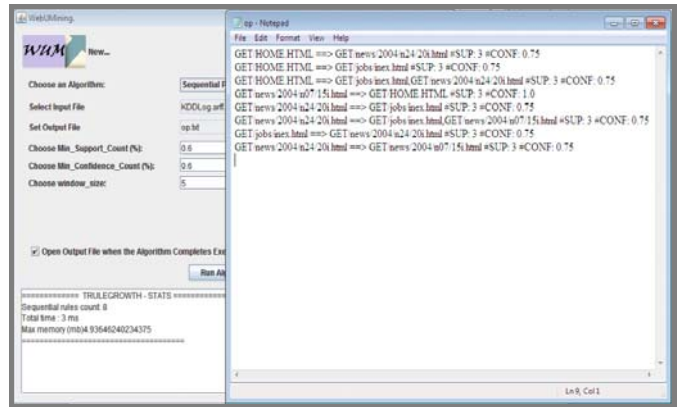
Figure 4: Result Screen of SequentialPattern algorithm


Figure 4: Result screen with different support value

## V. VISUALIZING SEQUENTIALPATTERN RESULT

The visualization is the vital part of the research work. It is used to compare the user access behaviour when processing web access log data through data mining techniques for achieving web usage mining process. When visualizing user access sequential pattern it is important to present the structure of the web site. Basically sequential pattern used to compare the start and end points that web user used during surfing the internet. It not only shows which pages accessed by users but also gives the order in which they are visited. Figure 5.6 shows the all sequential pattern that ate found in the sample data collected from the kdnugget's repository.
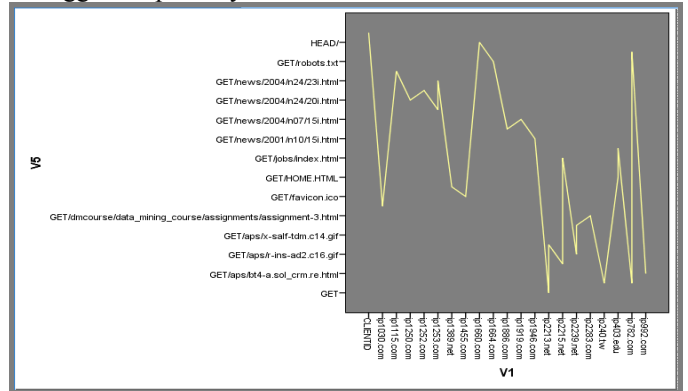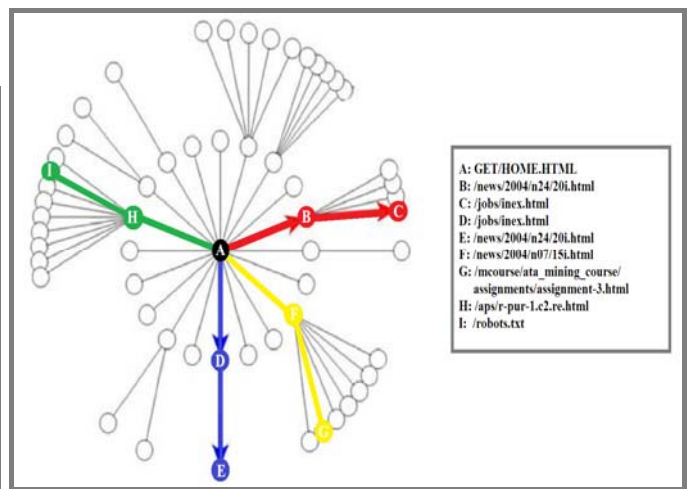

Figure 6: Visualization of Sequential Pattern


Figure 7: Comparison of Sequential Pattern

| No. of Sequential Pattern Generated | Time (in MSec.) | Minimum Support |
|---|---|---|
| 248 | 624 | 0.2 |
| 56 | 156 | 0.3 |
| 56 | 44 | 0.4 |
| 44 | 41 | 0.5 |
| 11 | 8 | 0.6 |

Table 3: Comparison of result generated by proposed SequentialPattern Algorithm

The table given above shows the sequential pattern generated by proposed SequentialPattern algorithm with respect to the minimum support measure and the time taken to generate the result so that it is useful to compare the sequential pattern with respect to the given support count.

## VI CONCLUSION

This research work proposed by the authors demonstrate a low level prototype called WebUMining which implentemented with the proposed SequentialPattern Algorithm which can be used to compare the user access behaviour using web access log data through web usage mining techniques. The proposed algorithm is implemented using Java and WEKA. For visualization of the result we have used the SPSS and RapidMiner software.

## REFERENCES

[1] Agrawal R. and Srikant R., "Mining sequential patterns", International Conference on Data Engineering , 1995.
[2] Yun U., "A New Framework for Detecting Weighted Sequential Patterns", Large Sequence Databases. Knowledge-Based Systems, 2008.
[3] Masseglia F., Teisseire M. and Poncelet P., "Sequential Pattern Mining: A Survey on Issues and Approaches", Encyclopedia of Data Warehousing and Mining", 2005.
[4] Gupta M. and Han J., "Approaches for Pattern Discovery Using Sequential Data Mining", Pattern Discovery using Sequential Data Mining, 2011.
[5] Tyagi N., Tyagi S. and Solanki A., "An Algorithmic approach to data preprocessing in Web usage mining", International Journal of Information Technology and Knowledge Management, 2010.
[6] Rathod K. and Valera," A Novel Approach of Mining Frequent Sequential Pattern from Customized Web Log Preprocessing", International Journal of Engineering Research and Applications, 2013
[7] Rajimol, "Data Mining For Web Intelligence", Thesis, MG University, 2013.
[8] Vasumathi, "Web Mining Using Pattern Discovery Techniques", Thesis, Jawaharlal Nehru Technological University, 2010.
[9] Han J. and Pet J., "Mining Frequent Patterns by Pattern-Growth: Methodology and Implications", SIGKDD Explorations, 2000.
[10] Cooley R., Mobasher B. and Srivastava J., "Web Mining: Information and Pattern Discovery on the World Wide Web", International Conference on Tool and Artificial Intelligence, 1997.
[11] Chen y. and Huang T., "Discovering Time-Interval Sequential Patterns in Sequence Databases", Expert Systems with Applications, 2003.
[12] http://www.kdnugget.com