# DOUBLE IMPLICIT WAPAPPROACH DESIGNED FOR SERIAL STRATEGY EXTRACTION FOR WEB ACTIVITY

[1]**Rahul Moriwal** and [2]**Dr. Anil Rao Pimpalapure**

[1,2]Eklavya University, Damoh (M.P.), India

## ABSTRACT

*With the exponential development of information accessible from the World Wide Web, discovering and analyzing meaningful data from online has become a pragmatic need. Online access patterns, or the series of reaches that clients commonly follow, are an intriguing and important piece of pragmatic insight. Successive Trend extraction is the practice of using data extracting techniques on an ordered record to identify correlation links between an ordered sequence of occurrences. Web access pattern tree (WAP-tree) extracting is a successive trend extraction approach for online activity reaches sequences. It begins by putting the initial online access order records on a prefix tree, same as to the frequent pattern tree (FP-tree) for unordered information.The WAP-tree algorithm then extracts the recurrent orders from the WAP-tree by iteratively creating intermediary trees, beginning with suffix series and ending with prefix series. A modification of the WAP tree technique was attempted in order to improve efficiency. The use of new Binary Coded WAP removes the demand for frequent re-constructions of intermediate WAP-trees during extracting, resulting in significantly reduced execution time.*

*Keywords: WAP tree, data extracting, successive information extraction, frequent pattern tree*

## I. INTRODUCTION

Data extraction is the complicated procedure of uncovering genuine, innovative, maybe beneficial and eventually understood trends in information. Organizations face challenge of information overload as records become more widely used and expand in size. The challenge of successfully utilizing these huge amounts of information is becoming an acute issue for all organizations.

Typically, we have used information to query an appropriate database repository through a well-defined utilization for prefabricated report generation usage. While this method of engagement is enough for a wide range of areas, exploratory analysis of information is needed for several other applications. These programs offer query-driven information consumption, which means that the analysis is based on a human analyst's question. On the other side, information extracting process allow for the automated investigation of information. Information extracting seeks to identify trends in data as well as deduce rules from them. The client will be able to utilize these guidelines to promote, assess, moreover evaluate judgments in a linked commercial field. This creates the prospect of a fresh approach to interface with records and information repositories.

Successive extraction is the practice of using information extraction algorithms on an ordered records to identify correlation links between a serial sequence of occurrences.

The purpose of this study is to use information extracting algorithm on a serial order to identify correlation links between an ordered list of occurrences. Given a WASD (Web Access Sequence Database), the purpose is to identify commonly occurring successive trends using the minimal amount of help available. Successive trend extracting is used in fields such as medical treatment, scientific and engineering operations, and telephone call patterns. Sequential pattern mining, online use mining is the automated finding of client access trends from online servers. It is used by online platforms to detect prospective clients who are prone to spend a lot of capital or to forecast which website visitors would click in certain ads or banners by observing how past visitors have reacted to such banners.

## II. BACKGROUND

Successive trend extraction is used in association rule extraction. The transaction database's association rule T is an expression of the type X Y, where X and Y are subsets of A. If the proportion of transactions in D that support X also supports Y, then X Y holds with trust. If a portion of the transactions in transaction set T support X U Y, then the rule X Y is supported in that transaction set. Mining association guidelines requires two steps. Initially, typical trends concerning the minimum help threshold (SUP) are mined. Secondly, the minimal trust and trust threshold are used to build association rules. There are two kinds of trend extracting:

**[1] Non-sequential trend Extracting** involves analyzing things in a transaction without regard for sequence.

**[2] Sequential Trend Extracting:** Items in a transaction follow a specific sequence and may repeat in the similar order.

The acronym WAP-tree stands for web access pattern tree. The major phases in this approach are outlined below. The WAP-tree integrates web activity information in a prefix tree manner, same as to the frequent pattern tree (FP-tree) for non-sequential information. The program analyzes the web activity once to identify common specific occurrences. Second, it does another search of the web log to build a WAP-tree over the grouping of recurring individual events for every transaction. Thirdly, it finds patterns of conditional suffixes. The intermediate conditional WAP-tree is constructed in the fourth step using the pattern from the preceding stage. Lastly, until the specified WAP-tree is either empty or contains a single branch, Steps 3 and 4 are repeated.

**Table I.** Sequence Database for Wap-Tree

| TID | Web access sequence | | Frequent Subsequence |
|---|---|---|---|
| 100 | pqspr | | pqpr |
| 200 | tptqrp | | pqrp |
| 300 | opqupt | | qpqp |
| 400 | puqprur | | pqprr |

Thus, using the WAP-tree approach, discovering all common occurrences in the web activity requires creating the WAP tree and extracting accessibility trends from it. The web activity ingress series order in Table 1 is used to demonstrate WAP-tree construction as well as mining. Assuming the lowest help criteria is set at 75%, which indicates that in our example, ingress series should contain a count of three out of four records to be considered frequent. In order to build a WAP-tree, the database must first be scanned once to gather common occurrences. Every sequence's irregular section is removed before building the WAP-tree. The input is limited to the subsequences that occur often. For example, Table 1 shows all events as p, q, r, s, t, and u, with p, q, r, s, t, and u having helping values of 4, q, s, t, and u. Only p, q, and r occur often when there is a lowest help of three. In order to create the recurrent subsequence shown in, all irregular events (such as s, t, and u) are eliminated from each transaction sequence illustrated in Column 3 of Table I.
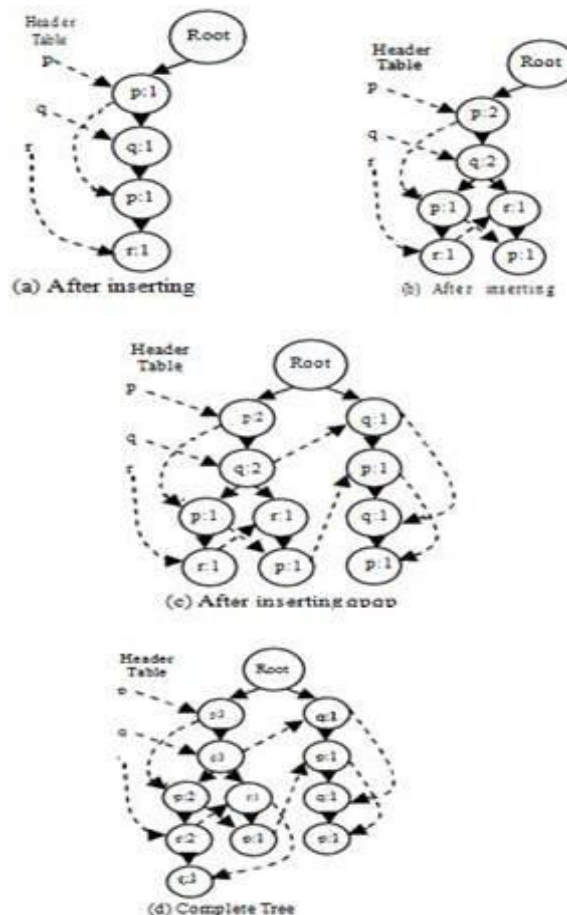
With the recurrent series in every operation, the WAP-tree method initially saves most common products as header nodes, which are then utilized to connect all nodes of their kind in the WAP-tree in the series they were added. When creating the WAP tree, a virtual root (Root) is initially added. Then, a trail from the Root to a leaf node of the tree is made using the transaction's frequent sequences. Every event in the sequence gets allocated a node with a count of 1 from Root if the node type is still unknown; if not, the node's count is increased by one. Additionally, the head link for the inserted event is connected (in broken lines) to the freshly added node from the previously added node of its type, or, if it is the first node of that event type added, from the header node of its kind. In figure 1(a), a

**Copyrights @ Roman Science Publications Ins.**                                    **Vol. 5 No.2, June, 2023**
**International Journal of Applied Engineering & Technology**

308

*International Journal of Applied Engineering & Technology*

left child of Root with label p and coun is constructed in order to add the first frequent sequence pqpr of transaction ID 100 in the sample database.

1.The added node from the p header is subsequently linked (in broken lines) to the header link node for regular occurrence p. The following event, q, is connected to header node q and placed as the left child of node p with a count of 1. With a count of 1, the third event, p, is positioned as the left child of node q, and the p link connects this node to the inserted node. The fourth and last event in this series is r, which has a count of one and a connection to the r header node. It is positioned as the left child of the second p on this branch. Second, start entering the subsequent transaction's series pqrp using ID 200 from the virtual Root (figure 1(b)). The count of node p is increased by one to obtain (p: 2) since the root includes a child called p. In a similar vein, the tree contains (q: 2). A new node, r:1, is created and inserted as an additional child of the q node because the next occurrence, r, is incompatible with the next existing node, p. Figures 1(c) and (d) are produced by inserting the third sequence (qpqp) of ID 300 and the fourth string (pqprr), accordingly.

Figures 1(c) and (d) are then created by adding the third sequence (qpqp) of ID 300 and the fourth sequence (pqprr).

Following the addition of the ordered information to the whole WAP-tree (figure 1(d)), the tree may be extracted for recurrent trends, starting with the header list's least regular event—in our case, frequent event r—as will be seen in the discussion that follows. The prefix series of base r or the required series base of c may be calculated using the WAP-tree in figure 1(d) in the following ways: pqp:2; pq:1; pqpr:1; pqp:-1.
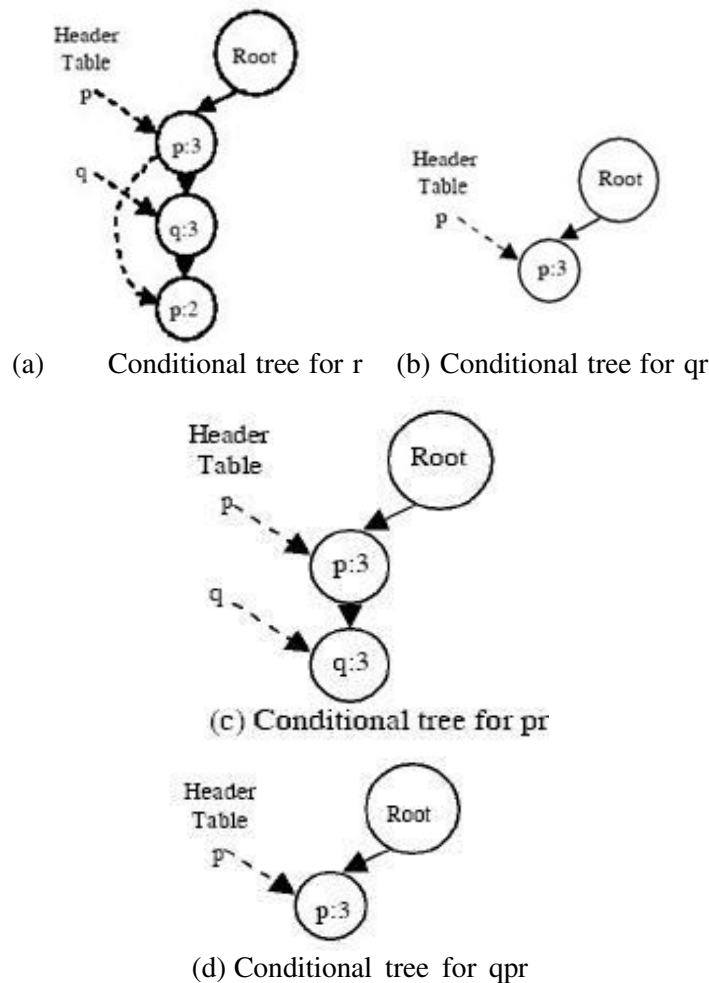


**Figure 1.** Construction of WAP Tree

---

## International Journal of Applied Engineering & Technology

WAP-tree|r is created in the similar way as illustrated in Figure 1. Figure 2(a) shows the new conditional WAP-tree. Recursively, using the WAP-tree in Figure 2(a), the

following required series base for the upcoming suffix sequence, qr, is discovered to be p(3). With p being the sole regular trend in this base, the recurrent series base of qr utilized to build the upcoming WAP tree as seen in figure 2(b) is p(3). This concludes the rebuilding of WAP trees that began as suffix series |r, |qr, and the most common trends identified along this line are r, qr, and pqr. The looping process continues using the suffix path |r,|pr. The suffix pr's required series basis is Ø, pq:3, as seen in figure 2(a). Figure 2(c) shows the WAP tree constructed from this list. The algorithm continues to run, determining the required series bases of qpr with p:3. From the list, the required regular occurrences of pqr are just p:3. Figure 2(d) shows how the conditional WAP-tree|qpr is created. Returning to the extracting of common trends with suffix pr, figure 2(c) mines for required series bases with suffix ppr and returns NULL.

The necessary inquiry for r is currently been finished. The extracting process for trends with suffix r is identical to that of searching for recurring trends with the suffix of other header regular occurrences (starting with suffix base |q and moving on to |p). Following tree extraction, the following trends set is most often encountered: {r, qpr, pqpr, pr, pqr, qr, qb, pq, p, pp, qp, pqp}.



(a)      Conditional tree for r     (b) Conditional tree for qr

(c) Conditional tree for pr

(d) Conditional tree for qpr

**Figure 2.** Reconstruction of WAP trees for mining conditional pattern base r

Copyrights @ Roman Science Publications Ins.                                  Vol. 5 No.2, June, 2023
International Journal of Applied Engineering & Technology

310

## *International Journal of Applied Engineering & Technology*

### III.  RELATED WORK

### A.  *Binary Coded Web Access  Trend  tree  Method*

Identical to a WAP-tree, the tree data structure is used to compactly store navigating records series and the associated counts of recurring events, hence preventing laborious help enumeration during extraction. Every node in the modified WAP-tree is given a binary code. These codes are employed in extraction to locate nodes inside the tree. The nodes are connected through a series of events to generate the header table. Here, linking is used to traverse prefix sequences while keeping track of nodes that have the same label. Instead of using suffix search, the extraction strategy uses prefix series finding.

### B.  *The Approach*

Input : Access series records D(i), limited help MS (0< MS d" 1)

Output : regular successive trends inD(i).

Variables : a saves whether a node is an ancestor in the queue,and Cn saves the entire amount of instances in the suffix trees.

Start

Scan D(i) to discover recurrent specific instances L; Scan D(i) repeatedly. Make a root node of Tree T. code(root)= NULL;

count = 0; {

For (every access series, fs in D(i)) {

Mine regular suborders ($fs_1 fs_2 \ldots fs_n$) by removing all instances that are not in L;

current node -> leftmost_Child(root);for ( k=1 to n ) {

if (current node = NULL)

{Make a fresh child node with position code equal to "1"appended

To position code of parent of current node ;}elseif (current node = $fs_k$ ) { NdFd =  true ;}

else { make current node point to current node sibling}

}

if (NdFd = true)

{count (fsk) ++;

Make current node point to $fs_k$ ;}

Else {make fresh child node with position code of current

node with

"0" appended at the end;

Make current node point to new created node ;} } }

From root node, do a successive Traversal of Tree T to create proper connect queue;

TREND_DIS (Suffix tree roots STR, Frequent sequenceFS);

end;

**Copyrights @ Roman Science Publications Ins.**                                    **Vol. 5 No.2, June, 2023**
**International Journal of Applied Engineering & Technology**

311

# *International Journal of Applied Engineering & Technology*

TREND_DIS(R, F) {

If (STR=empty) return;

for (each suffix tree of event in L) {Save first event in $e_i$      queue to A

if (event $e_i$     is descendent of any instance in STR, and is not descendent of A)

{Insert ei suffix tree header set STR';Add count of $e_i$      to Cn;

Replace the A with $e_{i.}$ ; }

If(Cn > MS )

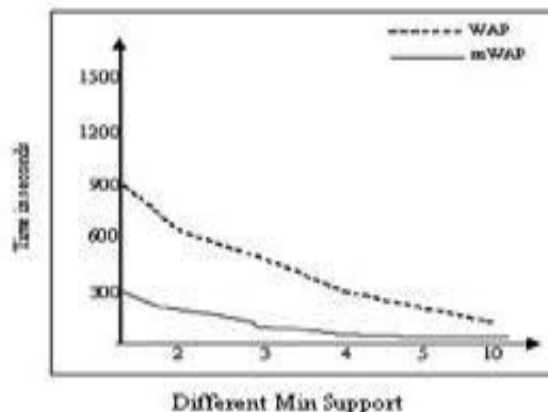{Append ei after FS to FS';print (FS');

TRENd_DIS (STR',FS'); }

## IV.EXPERIMENTAL RESULTS

This experiment makes use of a fixed-size record and varying limited help levels. The records as well as methods are evaluated against a database of 60 thousand (60 K) with minimal supports ranging from 0.8% to 10%.

Table 2 and Figure 7 show that: Table II. Execution Times For Records At Differentlimited Help

| Algorithms | time in secs at different supports | | | | |
|---|---|---|---|---|---|
| | 2 | 3 | 4 | 5 | 10 |
| WAP | 750 | 510 | 330 | 280 | 150 |
| Modified WAP | 230 | 160 | 110 | 95 | 48 |

As the lowest-highest help increases, each technique's execution time decreases. This is due to the fact that there are fewer possible series at less backing increases. Because of this, methods finding common patterns faster. The runtime of the modified WAP technique is consistently less than that of the WAP approach. Higher memory and input/output (I/O) storage costs result from WAP tree mining. The cost of maintaining intermediated trees considerably extends the program's entire execution time, even on systems with memory only. It makes more sense to assume that these methods are applied in common systems that are accessible from different places. These systems may not only be memory-based, but may also involve several processing systems that share memory and central processing units (CPUs) with the assistance of virtual memory.
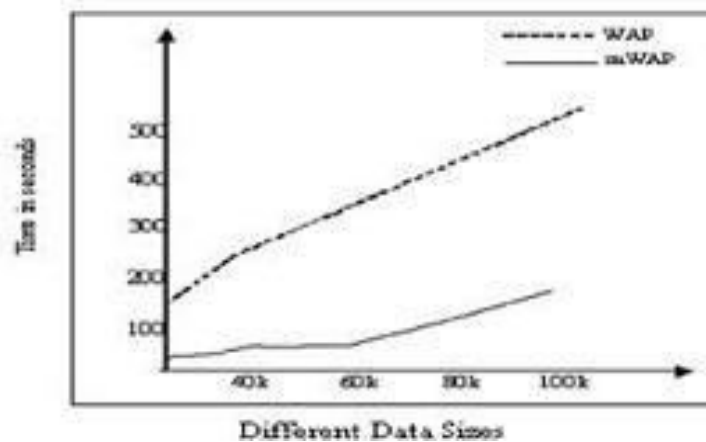


**Figure 3.** Execution times trend with different limited help

Copyrights @ Roman Science Publications Ins.                                                      Vol. 5 No.2, June, 2023
International Journal of Applied Engineering & Technology

312

# *International Journal of Applied Engineering & Technology*

Now, databases with different sizes from 20 K to 100 K with the fixed lowest help of 7% are used.

**Table III.** Execution Times Trend with Different Data Sizes

| Algorithms | Different changed transaction size | | | | |
|---|---|---|---|---|---|
| | 20k | 40k | 60k | 80k | 100k |
| time in sec | | | | | |
| WAP | 148 | 265 | 320 | 445 | 540 |
| Modified WAP | 50 | 75 | 97 | 145 | 179 |



**Figure 4.** Execution times trend with different data sizes

## V. CONCLUSION

In this paper, we look at the successive trend extraction field. The modified version clearly performs better than the web access trends tree strategy after comparing the two methods. In addition to a performance comparison with a graph, this section reviews the advantages and disadvantages of both approaches. With the new method, storing many intermediate WAP trees while extraction is no longer necessary. Particularly when mining very long series with millions of entries, it greatly decreases enormous memory access costs because only the original tree is preserved. This may require disk I/O in a virtual memory environment. By using this method, recreating interim WAP trees may be done without having to preserve and scan intermediate conditional pattern bases. In order to keep all events ei in the identical suffix tree close to one another in the linkage, this method uses pre-order linking of header nodes, which improves search efficiency. There has also developed a straightforward technique for assigning position codes to nodes in any tree, which can be utilized to ascertain the relationship between tree nodes without having to repeat the traverse.

## REFERENCES

[1] Pei, J., Han, J., Mortazavi-Asl, B., and Zhu, H., "Mining access patterns efficiently from web logs. In Proceedings of the Pacific-Asia Conference on Knowledge Discovery and Data Mining (PAKDD'00)", Tokyo, Japan, 2000.

[2] M. Eirinaki, M. Vazirgiann, is "SEWEP: Using Site Semantics and a Taxonomy to Enhance the Web Personalization Process", in Proceedings of the 9th ACM International Conference on Knowledge Discovery and Data Mining (SIGKDD'03), Washington DC, August 2003.

[3] W3C, "Logging control on W3C httpd.", www.w3.org/Daemon/User /Config/Logging.html #common-log file-format, July 1995.

Copyrights @ Roman Science Publications Ins.
Vol. 5 No.2, June, 2023
International Journal of Applied Engineering & Technology

313

## *International Journal of Applied Engineering & Technology*

[4]     R. Cooley, B. Mobasher, and J. Srivastava, "Web mining: Information and pattern discovery on the World Wide Web", In Proceedings of the 9th IEEE International Conference on Tools with Artificial Intelligence (ICTAI'97), 1997.

[5]     G. Salton, C. Buckley, "Term weighting approaches in automatic text retrieval", Information Processing and Management, Vol. 24, 1998, pages 513-523.

[6]     Robert Cooley, Bamshad Mobasher, Jaideep Srivastava, "Data preparation for mining world wide Web browsing patterns", Knowledge and Information Systems, Vol.1, No. 1, February 1999.

[7]     Han J. Pei, J. Yin, Mining frequent patterns without candidate generation. In Proc. ACMSIGMOD Int. Conf. Management of Data (SIGMOD'00), Dallas, TX, 2000, pp 1–12.

[8]     S. Taherizadeh, N. Moghadam, "Integrating Web Content Mining into web usage mining for finding patterns and predicting users behaviors", International Journal of Information Science and Management, Vol. 7, No.1, June 2009.

[9]     D. Mladenic, "Machine Learning used by Personal Web watcher", in Proc. Of ACAI-99 Workshop on Machine Learning and Intelligent Agents, China, Crete, July 1999.

[ 10]   Jose Borges, Mark Levene, "Data Mining of User Navigation Patterns, in Web Usage Analysis and User Profiling", published by Springer-Verlag as Lecture Notes in Computer Science, Vol. 1836, pages 92-111, 1999.

[11]    K. Wu, P.S. Yu, A. Ballman, "Speed Tracer: A Web usage mining and analysis tool", IBM Systems Journal, Vol. 37, No. 1, 1998.

[12]    Chen. T-S, Hsu. S, "Mining frequent tree-like patterns in large datasets", Data & Knowledge Engineering, pp. 65-83, 2007.

[13]    R. Agrawal and R. Srikant, "Mining Sequential Patterns power set", In ICDE, pages 3-14, 1995.

[14]    M. J. Zaki, "Parallel Sequence Mining on Shared-Memory Machines in Large-Scale Parallel Data Mining", M. J. Zaki and C.-T. Ho, editors, Lecture Notes in Artificial Intelligence (LNAI 1759), Springer-Verlag, Berlin, 2000.

[15]     Dong, G. and Li, J. 1999. Efficient mining of emerging patterns: Discovering trends and differences. In Proc. 1999 Int. Conf. Knowledge Discovery and Data Mining (KDD'99), San Diego, CA, pp. 43–52.

[16]    R. Srikant, and R. Agrawal, "Mining Sequential Patterns: Generalizations and Performance Improvements", In Fifth Int'l Conference on Extending Database Technology (EDBT'96), Avignon, France, March 1996, pages 3-17.

[17]    H. Mannila, H. Toivonen, and A. I. Verkamo, "Discovering Frequent Episodes in Sequences", In Proc. of the First Int'l Conference on Knowledge Discovery and Data Mining (KDD'95), Montreal, Quebec, pages 210-215, 1995.

[18]    Han, J., Pei, J., Yin, Y., and Mao, R., "Mining frequent patterns without candidate generation: A frequent pattern tree approach", International Journal of Data Mining and Knowledge Discovery. Kluwer Academic Publishers, 2004, pages 53–87.

[19]    F. Massegila, F. Cathala, and P. Poncelet, "The PSP Approach for Mining Sequential Patterns", In Proc. European Symposium on Principle of Data Mining and Knowledge Discovery (PKDD'98), Nantes, France, pages 176-184, September 1998.

**Copyrights @ Roman Science Publications Ins.**                                              **Vol. 5 No.2, June, 2023**
**International Journal of Applied Engineering & Technology**

314

[20] M. J. Zaki, "SPADE: An Efficient Algorithm for Mining Frequent Sequences", In Machine Learning, Vol. 40, pp. 31-60, 2001.

[21] J. Pei, J. Han, B. Mortazavi-Asl, H. Pinto, Q. Chen, U. Dayal, and M.-C. Hsu, "Prefix Span Mining Sequential Patterns Efficiently by Prefix-projected Pattern Growth", In 17th International conference of Data Engineering (ICDE'91), Heidelberg, Germany, Apr. 2001.

[22] Han, J., Pei, J., Mortazavi-Asl, B., Chen, Q., Dayal, U. and Hsu, "Free Span: Frequent Pattern-Projected Sequential Pattern Mining", Proc. Int'l Conf. on Knowledge Discovery and Data Mining (KDD'00), Boston, MA, 2000 pp.355-359.

[23] Ezeife, C. and Lu, Y. Pie, "Mining Web Log Sequential Patterns with Position Coded Preorder Linked WAP-Tree", International Journal of Data Mining and Knowledge Discovery (DMKD) Kluwer Publishers, pp.5-38, 2005.

[24] M. N. Garofalakis, R. Rastogi and K. Shim. "SPIRIT: Sequential Pattern Mining with Regular Expression Constraints". In VLDB, 1999, pages 223-234.

[25] M. Arnoux et al., "Automatic Clustering for the Web Usage Mining", Proc. 5th Int'l Workshop Symbolic and Numeric Algorithms for Scientific Computing (SYNASC 03), Editura Mirton, pages 54–66, 2003.

[26] Rene Arnulfo García-Hernandez, "Finding Maximal Sequential Patterns in Text Document Collections and Single Documents", Informatica 34 (2010) 93–101.

[27] HuaJiang, DanZuo, Xin Hu, Yong-XinGe, Bin Han, "UAP-Minar:A Real-Time Recommendation Algorithm Based on User Access Sequences", 7th Int'l Conf. on Machine Learning and Cybernetics, Kunming, 12-15 July 2008.

[28] Unil Yun and John J. Leggett, "WSpan: Weighted Sequential Pattern Mining in Large Sequence Databases", Proc. Of the Third Int'l Conf. on IEEE Intelligent Systems, Sep 2006.Pages 512-517.