

# Pattern Discovery of Usage Pattern from Web Data

Mr. Pratap Lal<sup>1</sup> Dr. Ajeet Kumar Singh<sup>2</sup>

<sup>1,2</sup>Department of Computer Science and Engineering

<sup>1,2</sup>Suyash Institute of Information Technology, Hakkabad Gorakhpur (UP), India

**Abstract**— Web usage mining is a data mining methods. There are huge amount of data are stored on the internet. When user find any particular information by search engine like Google, Bing etc. is very hard because the complexity of web pages is increases day by day. Web usage mining play a vital role to solve this problem. In web usage mining we are creating an appropriate pattern according to the user's visiting behavior. The aim of this dissertation report is to applying an association rule using FP-Growth algorithm on web server log file (an educational institution web log data) to find the behavioral pattern and profiles of users interacting with a web site. The web mining usage pattern of a Technical Institution web data. Web related data is coteries in to three parts namely web log, access log, and proxy log data and collect the data in web server and we analyze this data using web log analyzer tool and WEKA tool. Our experimental results help to predict and identify the number of visitor for the website and improve the website usability. This work using data from user visiting system web log files that produced from web server IIS (Internet Information-Services) and using the accessed web address page references and access time.

**Keywords:** Web Mining, Web Access Logs, Web Usage Mining, Access Log Analyzer, World Wide Web, Weblog Expert

## I. INTRODUCTION

The Internet plays an essential role to hold, share and distribute the information. Due to quick development of the web pages is give a great chance to find the user and system behavior due to learn web access logs. Can access the document, modify or delete whole document according to privilege in real time.

Data Mining is the method to determine and search to unknown content automatically from the web server. Web is not only used for a strong medium to searching helpful information but also it used for searching information from web server. By the definition of web mining there are two methods can be used namely presses based and data based.

### A. Web Usage Mining:

Web uses mining plays an important role for extracting helpful information and finding suitable pattern. These patterns are more useful when user search any particular web based application. WUM is some time known as web log mining. It has three integrated part namely preprocessing, pattern discovery, and pattern analysis [1]. These three parts are take input data (log server data) one by one.

### B. Process of Web Usage Mining:

Web-usage mining is subjected to data mining process and it restive the previously unknown information from web server. Web Usage Mining plays an effective role to improve availability web sites. The web usage mining contain three different phases to perform process of web mining (1)Pre-

processing phase (2) pattern discovery phase, (3) pattern analysis phase.

### C. Application of Web Usage Mining:

The quantity of WUM applications is developing persistently, because of the business enthusiasm for online business Web locales and the related Web-promoting applications. Also, the developing enthusiasm for Web semantic field and the current field of Web semantic mining will bring new viewpoints for the WUM-related applications (Facca and Lanzi 2005, Srivatava et al 2000).

## II. LITERATURE REVIEW

### A. Web Usage Mining:

WUM use data mining technique to extract the useful information web server database according to the user request. There are mainly three steps are used namely preprocessing pattern discovery and pattern analysis phase. The main aim of web usage mining is to develop the whole configuration of website, to progress the availability of info material of the site.

According to [12] WUM is a very interesting sector for the research. Web-mining is the technique or part of data mining that solves the web-site performance problem and optimize the availability of information.it can used some data mining algorithm and association rule to fine the pattern according to user behavior. They also define the web usage mining is categorize into 3 large category such as web-content mining, web-structure mining and web-usage mining. C.C. Chan (2006) is also describe the WUM is subjected to data mining process. That fined previously unknown information automatically from the web server and also define pattern discovery rule to improve the web site performance.

### B. Pattern Discovery:

As per D. Jagli and S. Oswal,(2012) [26] was portray the example revelation process is an information mining system and furthermore clarify the procedure of web use mining. The web use mining is isolated into three in number classifications such a web content mining, web structure mining and web-usage mining. The WUM is an exceptionally viable fled of information mining it additionally have three levels working procedure. Pre-processing Pattern Discovery and Pattern Analysis.

S.Prakash and R.M.S. Parvathi [9] Discussed a successful arrangement of standards utilizing affiliation manage (convention) and Apriori calculation. This convention suit powerful to make example and utilizing this example run improved the accessibility of obscure data. The new mining principle is same as the set up example affiliation and Apriori calculation.

Kumar et. al [10] concentrated on finding of log information at a web server on find the utilization examples of site from log documents. They executes vital three

countenances of web use mining apriori calculation create affiliation decide that bury relates the customer's utilization design for specific site. They actualizes FP-Growth calculation. Apriori is actualize to work the database contains set off exchange.

Agrawal et al. [5] suggested Apriore algorithm.it is an effective data generating algorithm that are used during candidates are generated. It is based BFS algorithm model to count item sets in data generation module.

### III. PROBLEM DESCRIPTION

Web usage mining is playing an important role to solve this problem. There are lots number of research work are also done in this filed. WUM processes have three stages i.e. Pre-processing, Pattern-discovery and pattern-analysis. Our work is based on Pattern discovery phase.

My main aim is to making a system that printout (identify) user visiting pattern from weblog data of a college web site. There are frequent candidate generating algorithms namely Apriori and FP-growth are used for frequent candidate generation that is sufficient for this work.

Apriori algorithm has some limits like it scan large database, high cost for large amount of frequent candidate generation, large amount of comparison (matching) from candidate sets. So due to this region Apriori algorithm is more costly time taking and inefficient for this work. So in order to solution of this problem is to FP-growth algorithm is better.

FP-growth algorithm has two steps, in first step is construct FP-tree and in second step search FP-tree and give output for all deferent patterns.

Explosive lack is main disadvantage of FP growth algorithm that is use for respectable candidate generation technique.

### IV. METHODOLOGY

According to Jaideep, et al (2000), web usage mining has three main processes to reveal knowledge out of the data warehouse or log file (see figure 4). As discussed in chapter two, there are various tools and algorithms for mining web logs.

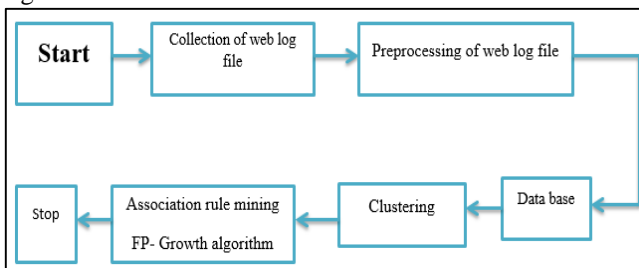


Fig. 1: Proposed Framework

#### A. Description of the Model

##### 1) Web Log File:

The row web log data is collected from weblog server. Weblog server store every activity of user that visits on web site and this web log record is used to analyze user visiting behavior on web site and according to this record we improve web site design .web log data stored many different formats like extended log and mutual log format. A sample of web log

data is shown in figure 4.2 and this data is collected from the college Web server.

Start Date	2017-02-08 00:00:00				
End Date	2017-02-08 23:59:59				
View Server Time	Wed Feb 08 15:25:33 GMT+05:30 2017				
Appliance	CR2500iNG-XP				
Firmware Version	10.6.4MR-1				
Appliance Key	C47316190997-SH2CRX				
Time	User Name	User Group	DomainURL	Category	IP Address
017-02-08 15:21:47	ec2015041028	student15-16	www.chiplensvctennis.net	None	172.20.1.165
017-02-08 15:21:47	raj	Student	*googleusercontent.com	*googleusercontent.com	InformationTechnology 172.16.3.40

##### 2) Data pre-processing:

Information pre-processing contain three most imperative strides specifically information cleaning, client recognizable proof and session ID. In pre-processing we expel superfluous sections, copy and repetitive dishes like us graphical and sight and sound passages .client and session distinguishing proof is an imperative assignment of pre-processing stage in this stage we evacuate the copy IP addresses and discover a gathering of comparative session applicants.

##### B. Pattern Discovery:

Pattern discovery is deals with second stage of Web Usage Mining and this is applied after pre-processing. The main task of pattern discovery is to finding pattern generating rule by using some techniques like association and clustering. The problem of association is described by agrawal [3]

##### 1) Association Rule Mining:

Affiliation manage mining issue was indicated by Agrawal [3]. Affiliation administer mining is one of the information mining procedures which is utilized to find valuable example. It takes a shot at producing successive example and tenets. In web log document number of URL visit by number of clients so we can recognize much of the time got to site pages by clients which can comprehend client needs. Two essential parameters of affiliation manage are support and certainty. We can characterize affiliation manage as takes after:

Let  $I = \{i_1, i_2, \dots, i_m\}$  be an arrangement of literals or things,  $D = \{t_1, t_2, \dots, t_n\}$  be an arrangement of exchanges, where every exchange  $t_i$  is a thing set with the end goal that  $t_i \in I$ . Every exchange,  $t$ , has an exchange id ( $t.id$ ) and a thing set ( $t.Itemset$ ), i.e.,  $t = (t.id, t.Itemset)$ . A exchange  $t$  contains an itemset  $X$  if  $X$  is a subset of  $t.Itemset$ . An Association run,  $R$ , signified by  $R: X \rightarrow Y$ , where  $X$  and  $Y$  are itemsets that wear  $t$  cross. Each manage  $R$  has two esteem measures, support and certainty, meant by  $sup(R)$  and  $conf(R)$  individually. The support of a thing set,  $X$ , has bolster,  $s$ , in exchange set,  $D$ , if  $s\%$  of exchange in  $D$  contain  $X$ .

Then, support is  $(R : X \rightarrow Y) = support(X \cup Y)$ , confidence  $(R : X, Y) = support(X \text{ and } Y) / support(X)$

Certain transactions contains similar item-set, particularly for remote-sensed images. This advises a approach to remove duplicate calculations. Certain ideas are given. Let  $T = \{n \mid n \text{ is some probable transactions}\}$ , while  $M = \{n \mid n \text{ is a transaction previously occurred}\}$ .

Input: Transaction gives the combination of all transaction

Output: every frequently item group can contain lesser threshold status

Begin

```

Allow numbers to every item
For each item a, Convert transaction table into bit vector
For each item a
{
If (number of 1's achieve minimum support condition of item
a)
{
Add this item in large 1-itemsets
}
}
Generate candidate 2-item-sets using large 1- item-sets
For each candidate 2-item-sets (a, b)
{
If (number of 1's is achieve minimum support condition after
logical AND operation of item a and b)
{
Add it to large 2-item-sets and create an edge between vertex
a to b
}
}
For each large n-item-sets (a1; a2; . . . ; an)
    Procedure (algorithm) for Association-Rule mining
    
```

**C. FP-Growth Algorithm:**

The FP-Growth algorithm produces frequent data sets from FP Tree by navigating in a bottom up approach [6]. This method decreases the total number of user data sets by generating a compacted type of database in terms of an FP-tree. This frequent information allows for the effective discovery of frequent data sets. It is a two-step approach and faster than other association mining algorithms [1].

**1) Algorithm (FP-Tree structure)**

Input: An improved entrée database and a smallest threshold minsup[FP].

Output: frequent pattern tree, FP-Tree.

Method: We applied the following steps to build FP-Tree.

- 1) Gather the combination of frequent articles F and supports. To Sort F in descending order as L, the list of frequent items. Test the access database once.
- 2) Firstly we build the root of FP-tree, T, and make it as "null". For each access data in table(a) in database do the following.

Select and sort the frequent objects in according to order of Length of tree L. Let the sorted frequent item list in table(a) be [p|P], where p is the first element and P is the enduring list. Call insert\_tree([p|P], T). The function insert\_tree([p|P], T) is executed as follows. If T has a child N so that N.item-name = p.item-name, then increment N's count by 1; otherwise create a new node N and let its count be 1, its parent link be linked to T, and its node-link be linked to the nodes with the same. The algorithm for constructing the FP-tree (Han et al., 2000) Example: It is a transaction IP address URL accessed and time database it use minimum-support threshold and calculate all frequent candidates.

IP Address	URL Accessed	Time	
T1 res1.ne.wi.ac.uk	/api/java.io.BufferedWriter.html	A	01/Jan/2017
	/api/java.util.zip.CRC32.html	B	01/Jan/2017
	/api/java.io.BufferedWriter.html	A	02/Feb/2017
	/java-tutorial/w/animLoop.html	C	04/Feb/2017
	/atm/logiciels.html	D	18/Feb/2017
T2 Acasun.cekerd.edu	/perl/pentre.html	F	11/Jan/2017
	/java/tutorial/animLoop.html	C	12/Jan/2017
	/html4_0/struct/global.html	G	29/Jan/2017
	/api/java.util.zip.CRC32.html	B	29/Jan/2017
	/postgres/html-manual/query.html	H	29/Jan/2017
T3 Acccs.francomedia.gc.ca	/java/tutorial/animLoop.html	C	05/Jan/2017
	/apache/manual/misc/API.html	I	05/Jan/2017
	/postgres/html-manual/query.html	H	05/Jan/2017
	/perl/pentre.html	F	12/Feb/2017
	/api/java.io.BufferedWriter.html	A	12/Feb/2017
T4 ach3.pharma.mcgill.ca	/api/java.io.BufferedWriter.html	A	05/Jan/2017
	/java-tutorial/w/animLoop.html	C	05/Jan/2017
	/html4_0/struct/global.html	G	05/Jan/2017
	/postgres/html-manual/query.html	H	12/Feb/2017
	/relnotes/deprecatedlist.html	E	12/Feb/2017

Table 4.3.1: An example of user access database

Transaction	Item Set
T1	A,B,C,D,E
T2	F,C,G,B,H
T3	C,I,H,F,A
T4	A,C,G,H,E

Table 4.3.1: (a). Coded form of user access database

Item	Support
A	3
B	2
C	4
D	1
E	2
F	2
G	2
H	3
I	1

Table 4.3.1: (b)

L1= The item set which are frequently repeated using minimum count.

Frequently Repeated Item Set	Support
C	4
A	3
H	3
B	2
E	2
F	2
G	2

Table 4.1: (c).

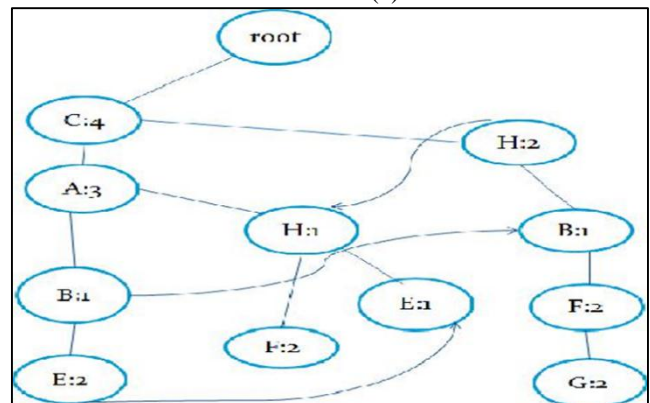


Fig. 2: The process of FP-Tree algorithm.

So finally we create pattern

- T1 : C->A->B->E
- T2 : C->H->B->F
- T3 : C->A->H->F
- T4 : C->A->H->E->G

### V. EXPERIMENT & RESULT

In this section we provide facts of the results obtained for each experiment. The experimental result is produce WEKA-tool .it is an open source tool for web mining. It take web log file in csv, and arff. Format and perform the operations of we usage mining like pre-processing, classification, association. In association phase it perform Apriori and FP-Growth algorithm to generate frequent sequential candidates.

#### A. Clustering:

The association-rules produced by the FP-Growth algorithm were analyzed to discover their interesting pattern tree

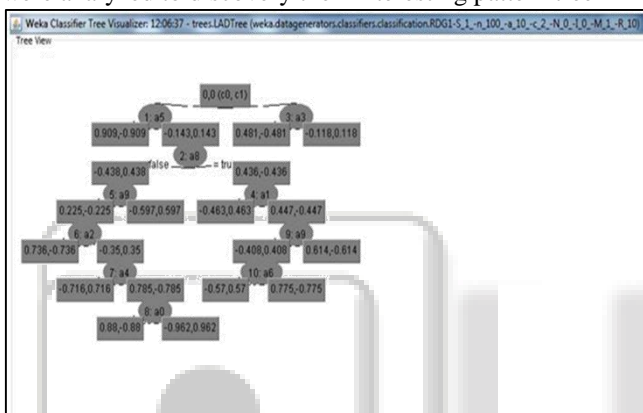


Fig. 2: Decision Tree of the 20-Most Frequent Time Features

#### B. Association Rules:

Association rules produced by the Apriori-algorithm hear it analyse and to find their interesting patterns

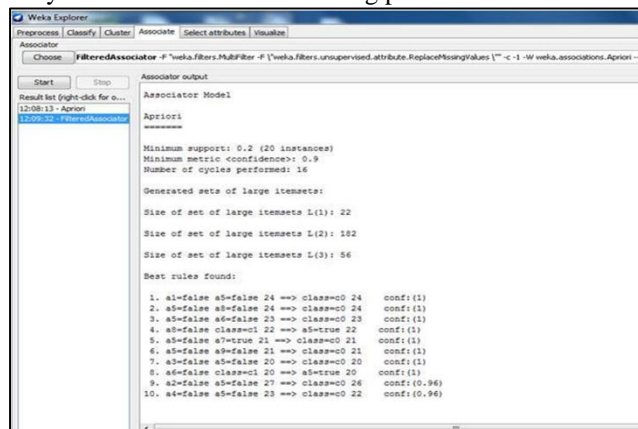


Fig. 3: Partial Output of the WEKA Apriori Program.

### VI. CONCLUSION & FUTURE WORK

This research work proposes a web usage mining frame for automatic web user access pattern based on frequent sequential pattern mining techniques in web log. This framework is capable of extract web user access pattern from large data sets with a very low support. The main objective of

this thesis is to find frequent-sequential pattern of a college website from server weblog data according to frequent user access pattern. The FP-Growth algorithm used to generate frequent-sequential pattern. FP-Growth algorithm is give better result comparison to Apriori algorithm. This frame work is the coherence that should exist among its 3 phases: 1-preprocessing phases, 2-pattern discovery phases and 3-pattern analysis phases. The three steps should be considered as being part of one single process and therefore, when defining the pre-processing options, it take maximum time of web mining process.it is also very difficult because it face data cleaning, user identification, and session identification. and in pattern discovery phase we are using FP-Growth algorithm witch is better than Apriori algorithm. Because it take less time space comparison to Ariori algorithm.

The proposed FP-Growth algorithm is used to improve the efficiency and accuracy of pattern. This is based on only support count rules, a chain of pre-processing works is conducted and the pre-processing data is converted to data cleaning and clustering format. Web access frequent sequence pattern mining is used to decrease interfering, analyzing the mining results and propose valued ideas on the improvement of the website.

For mining successive examples from web logs, the accompanying angles might be considered for future work. The method for changing the web log to successive database is still tedious and could be enhanced for web log mining. The Proposed calculation could be stretched out to deal with successive example mining in huge customary databases other than web log. For mining consecutive examples in exchange databases, there is a need to deal with simultaneousness of occasions.

Effective web use mining could profit by relating use to the substance of pages. Different zones of enthusiasm for future work incorporate dispersed mining with proposed strategy and applying these procedures to incremental mining of web logs and successive examples. Future research may concentrate on Web log data that incorporated with web substance and web linkage structure mining. This to help site page positioning, web archive order and the development of web data base.

### REFERENCES

- [1] Sanjeev, and Swati Goel. "Web Usage Mining: Finding Usage Patterns from Web Logs." American International Journal of Research in Science, Technology, Engineering & Mathematics (2013): 203-207.
- [2] Shukla, Rajesh, Sanjay Silakari, and P. K. Chande. "Web Personalization Systems and Web Usage Mining: A Review." International Journal of Computer Applications 72, no. 21 (2013).
- [3] Nina, Shahnaz Parvin, Mahmudur Rahman, Khairul Islam Bhuiyan, and Khandakar Entenam Unayes Ahmed. "Pattern discovery of web usage mining." In Computer Technology and Development, 2009. ICCTD'09. International Conference on, vol. 1, pp. 499-503. IEEE, 2009.
- [4] Suneetha, K. R., and Raghuraman Krishnamoorthi. "Identifying user behavior by analyzing web server access log file." IJCSNS International Journal of

- Computer Science and Network Security 9, no. 4 (2009): 327-332.
- [5] Baoyao, Zhou. "Intelligent Web Usage Mining." Nanyang Technological University, Division of Information Systems, School of Computer Engineering 94 (2004).
- [6] Santra, A. K., and S. Jayasudha. "Classification of web log data to identify interested users using Naïve Bayesian classification." *International Journal of Computer Science Issues* 9, no. 1 (2012): 381-387.
- [7] Facca, Federico Michele, and Pier Luca Lanzi. "Recent developments in web usage mining research." In *International Conference on Data Warehousing and Knowledge Discovery*, pp. 140-150. Springer Berlin Heidelberg, 2003.
- [8] Turban, Efraim, Ramesh Sharda, Jay E. Aronson, and David King. *Business Intelligence: um enfoque gerencial para a inteligência do negócio*. Bookman Editora, 2009.
- [9] Tyagi, Navin Kumar, A. K. Solanki, and Sanjay Tyagi. "An algorithmic approach to data preprocessing in web usage mining." *International journal of information technology and knowledge management* 2, no. 2 (2010): 279-283.
- [10] Upadhyay, Akshay, and Balram Purswani. "Web usage mining has pattern discovery." *International Journal of Scientific and Research Publications* 3, no. 2 (2013): 1-4.
- [11] Mr. Rahul Mishra, Ms. Abha Choubey, "Discovery of Frequent Patterns from Web Log Data by using FP-Growth algorithm for Web Usage Mining", *International Journal of Advanced Research in Computer Science and Software Engineering* || Volume 2, Issue 9, September 2012 || ISSN: 2277 128X
- [12] M. Parekh, A. S. Patel, S. J. Parmar, Prof. V. R. Patel "Web usage mining: frequent pattern generation using association rule mining and clustering", *International Journal of Engineering Research & Technology*, ISSN: 2278-0181, Vol 4, Issue 04, pp. 1243-1246, April 2015