



AENSI Journals

Australian Journal of Basic and Applied Sciences

ISSN:1991-8178

Journal home page: www.ajbasweb.com



Growth and Performance monitoring of Web-log file using CBFP and generating Cache hits for query templates

¹G.M. Karthik and ²Dr.S.Karthik

*1*Department of Computer Science and Engineering, SACS MAVMM Engineering College, Madurai, Tamil Nadu, INDIA
e-mail: gmkarthik16@gmail.com, Tel +91-94432-85673, +91-452-2389873

*2*Professor & Dean, Department of Computer Science and Engineering, SNS College of Technology, Coimbatore, Tamil Nadu, India

ARTICLE INFO

Article history:

Received 19 September 2014

Received in revised form

19 November 2014

Accepted 22 December 2014

Available online 2 January 2015

Keywords:

Web logs, Consensus tree, CBFP, Cache hits and usefulness, Template generation

ABSTRACT

The W3C provides access to a huge volume of data and the entire data present in it can be either identified or searched with the help of a search engine. The role of a search engine (used for retrieving the data stored in it) is unanimous. The server supported with the log file enables access to the data that is to be searched and accessed earlier. The role of FP tree provides access to the data along with user information and finally creates a web log file. In this effort the studies was focused with the pattern search which is a means to access enormous data volume and provide effective search result in the form of FP tree. The effectiveness of this study is the template creation, which was created then and there and hence the growth of the FP tree was not affected. The result of the study provides lower complexity irrespective of the data size that is provided as input parameter. The updated result was tested with enormous data out of which a minimum of 1, 00,000 test data provided effective result on analysis

© 2015 AENSI Publisher All rights reserved.

To Cite This Article: G.M. Karthik and Dr.S.Karthik., Growth and Performance monitoring of Web-log file using CBFP and generating Cache hits for query templates. *Aust. J. Basic & Appl. Sci.*, 9(1): 194-204, 2015

INTRODUCTION

The scenario of data usage and retrieval by users who are accessing W3C for information sharing and accessing makes the downpour of enormous data. The data base size differs based on users search and retrieval during the course of action. The unformatted data received has to be segregated with data filter. In order to provide interface Etzioni.O(1994) between the data and the user, despite the functionality of search engine it needs high data centric with various concentration proportion and this concentration can be extended to search engine automated answering system Scheffer.T (2004). Irrespective of the web structured data, most of the data search engine mechanism works with data mining concept searching data using web mining Senthil Pandian (2014) for frequent item Agarwal.R (2001) and information search.

The search engine process come across few hitches like overabundance in which data volume is really huge and it is difficult to trace the exact data due to complexity of data volume comes with valid and invalid data (subsets) sets to limited coverage. Complexity grows as data volume grows. The searching process was handled by query and the purpose is to provide index for direct accessing to the data and it leads to customized search with the tag name limited query followed by the user who posted the query.

Generalize approach of web search engine Cooley.R(1999) allows the search engine to be customized by the user with the keyword to be searched in the search engine and result in hierarchical tree which consists of several levels of various data, articles to be searched. The queries posted to categorize equivocal documents are very short to possess sufficient information. Such approach is cross checked with server side web log files along with client machine IP, page origination, time of request and the agent information for updating browser activities for user query in particular. The web log file grows as the search continues based on time. This followed by pattern utilization for performance monitoring. For producing better result we have addressed this using FP tree data structure for providing frequent pattern on the data stored in the web log files. While citing the work of Grahne.G (2000), the usage of web templates TRIE mechanism from web log file was well addressed in the work of Han(1999). Using prefix tree stores relevant quantitative information of user web log interface classified using consensus tree for rapid access of data's in templates.

Corresponding Author: G.M. Karthik, Department of Computer Science and Engineering, SACS MAVMM Engineering College, Madurai, Tamil Nadu, INDIA
Tel: +91-94432-85673; E-mail: gmkarthik16@gmail.com,

The nodes of the tree possess the Information of pattern and tree traversal provides other information of the page like keyword, search time, search through put and other valid information for seen with index buckets to store frequent pattern and its mutations for better selectivity.

The work of DeRaedt(2001),Deerwester.S (1990),Han(1993), Garofalakis(1999) andBerners-Lee (1994) projects the impact of CBFP will negotiate the growth of FP tree and some of the concepts like TRIE that deals with the principles of consensus tree followed by rule based constraint whose are the work of Han(1999,2000)respectively. This is for applying mutations and classification as addressed in the work of Thirumoorthy (2014) where they have addressed with k-nearest neighbor classification[46]for accessing the leaf node of FP tree. Incorporating all the concepts evolves our work using FP tree using constraint based frequent pattern mining accessing template of dynamic web log file.

The results of CBFP are the web template that answers the queries. The queries may be of composite of generalized keyword, data relevance based on keyword, deploying narrow based search with respect to time and the complexity of the search query. The format of CBFP is the resultant of TRIE structure for accessing maximum keyword stored in the FP tree (generalized query) to speed up the process. In addition CBFP provides clear visibility of search template files shown in the format of FP tree that address dynamic search process made with new search in comparison with consensus tree. This adoption of search strategy is based on particular sequence search progress which is dynamic based on user interest.

It starts with input data with frequent data matching followed by the instance in the search template. The template holds the information for all accessible keyword with certainty of data found in the web log file. It also deals with CBFP for defining user constraint by searching the desired keyword in minimum time.

The contemporary issues of keyword search in the web log file to access the data frequency was done with frequent pattern data structure by Craven(1999). This was done with two major processes like browsing and searching to add up with updating data structure for easiest and fast retrieval of data present in the weblog file. Ample work was done; in this the counterpart addresses the importance of updating data structures.

Updating web log file will never store previous updates and its lead to data misconception. User accessibility differs for one another so as the web log file. The updated web log file uses CBFP (proposed) for data search in templates with its frequency updates.

The organizations of this work were addressed by novel data mining method, for finding the facts and figure of the updates in the web log file. The resultant web log files stores the log information of present and past and were displayed in the form of FP data structures. The data traversal path for searching a data takes less time by neglecting unwanted data and its web templates. The proposed work significantly improves the weblog performance that satisfies user data and its updates.

On literature survey, CBFP algorithm with TRIE data structure, FP tree growth process, identifying the data in weblog file displays the search result in FP tree and finally with future enhancement.

Literature Survey:

The existing work addresses the concept of content and structure mining David (2008),Balabanovic.M (1995). The implication of web log data structure is dealt in Chen Li (2008). The links and its prediction using Web watcher systems Armstrong.R (1995) its user prediction in the form of link referred interested and non interested pages are focused in Armstrong (1995). Similarly the predicted or interested user link concept by David (2008) was addressed in the proposed work. Work on web utilization miner work is based on association rule that provides weblog file using advanced search to grasp the user preference identified by users' clicks made was addressed with user click streams concept Cooley (1999).The result of this work provides the user browsing path and analyzes the information based on user's browsing. The sharing of information using mining and analyze user preferences will prevent from breaches. Using association rule mining the webpage is reconstructed based on user queries Etzioni.O(1994), conceptual analysis for analyzing browser behavior using personalized web search Cooley (1999, 2000) displaying the content based on user interest with semantics Chen Li (2008), clustering and classifying weblog data Haibin Liu (2007), frequent pattern mining, search engine optimization Boyan.J (1996), mining challenges for easier browsing Craven (1999).

The double prediction by partial match scheme Heung Ki Lee (2009) using Apache provides efficient memory aware system. The concept is linked list tree for web accessed data was available with active patterns and done with recursive algorithm. The idea of web traversal pattern mining analyzes past and present using re-mining concept Jia-Ching.Y (2009) and the result was achieved with link matrix. Overcoming of supervised learning Nguyen.M.T (2010) removes the problem of handle sequence based on selective link attributes. The large volumes of data stored in the web log file Hawwash.B (2010) affected by un-supervised learning, were scalability and profile tracking. The inference was Rui.W (2010)Panagiotis.G(2010) in web log it does not provide efficient result in case of enormous data.

The need for synthetic data generator uses Dependency Graph Griffioen.J (1995) and Prediction for compression by partial match and provides result for predictive pre-fetching. The Dependency graph

performance becomes low and its accuracy is less when there is increase in the web log file size and its due fixed length transaction for accessing data from the web log file based on user constraint.

With the above insight the focus is to identify the CBFP based template and retrieve cache hits irrespective of size. The resultant template would be efficient and accurate providing exact information based on constraint. The constraint level focus on user interest and the data to be discovered and displays the data in FP tree format. In extent, CBFP patterns with data structure, the user constraint, and their update improves server performance.

Trie Structure:

TRIE projects the data retrieved from web log file in the format retrieved with dependencies of user interest with updated CBFP pattern. This data is generated automatically by the web server with its operations. To investigate on server performance, the pattern and its usage of frequent pattern mining (FP) was implemented. FP tree growth projects on TRIE structure with user interest (Constraint based).

Activity log file entry:

Log file entry targets the log knowledge which preprocesses the data gathered with the help of mining algorithm. Generally it accepts 3 kinds of input parameter as (i) Frequent pattern – which accepts frequently visited pages (ii) its time interval (iii) finally the mining sequence involved. The result of these parameters creates frequent pattern template that focus on various topics the user inquisitive about with different dimension showing alternate ways. Previous work produces result till pattern analysis and our work produces result of Frequent pattern Generations with template generation shown in Fig 1. The web mining program consists of quest engine with answers for every queries raised to the server as addressed by Senthil Pandian (2014). The IIS log file entries produces results based on user access to the server. Every entry has session symbol, timestamp and the URL of the visited web content.

User trails will be reconstructed and mistreatment of the methodology as Outlined by White and Drucker(2007).Here a tendency to extract arises from the logs of windows toolbar users and then user navigation is summarized with the help of browser path. The search trail that originates of question submission in search engine is measured by trails square. Hence, these search trails become the main source to coach the algorithms in case a similar question is repeated which is dealt in the following section. Every entry has session symbol, timestamp and the url of the visited web content. The session is balanced with session id, session time, and the domain URL.

1) Re-submission of question in search engine as discussed by Bilenko and White (2008) until the search engine terminates the process until it displays the information satisfying the search criteria.

2) The hyperlinks sequence is either checked for its consistency or multi-level search till the desire search strategy is met. The strategy deals with search path by eliminating the quantity of noise [Generally, a noise is a different paths that arises while searching for a relevant task or keyword by a user].

Pre-processing of data involves five steps per search path which are; log entry (associate anonymous session symbol, a timestamp and URL). The IIS records user activities data based on search processes that is supported by pre-processing for generalizing data mining formula. The net log provides the breakup knowledge of protracted sequence of users URL sessions. the retrieval of associated user session is arrived by sequential netlog entries(which will be two).In data-cleaning, users could share an equivalent IP address and also the same user can also appoint a totally different IP address at different times therefore it leads to a completely different sequence of visit at varied times. Lastly it satisfies the requirement of knowledge cleanup. Most add mining use a predefined interval to seek out the visiting session and it is attainable to get a lot of data concerning user sessions than employing a fastened interval. The learning of website grouping leads to bunch analysis which is highlighted by Lou *et al.*(2002) and the ways of using the same under cluster analysis. The session wise user and webpage updates are known when user switches his session. Finally, mistreating this data to chop and choose the net pages, a lot of correct user session data is obtained.

To identify the user: spotting the user by themselves and also through the cookies is the way in which the identification happens among the user and guest. Apart from that the domain id also contributes in identification. We've enforced a lot of general technique to spot user supported as shown by Cooley (1999). The follow up is that the last criterion that insists of FP tree origination with user identifications with the visited page.

Working with Session Management:

The user session has relationship with the Pattern discovery result because any identical user at the intervals the length of one specific website visit corresponds to user session. In line with this the citations Han (2000), Perkowitz.M (2000) has focused on page square measure. This measure should not exceed the pages visited stored with the session. We tend to consider the log data as a sequence as discussed by Ou *et.al* (2005) of distinct sites, where subsequent user sessions are discovered by remarkably long gaps between consecutive requests.

The Table 1 reveals that prime focus is on technical problems and is frail structured of weblog file information. The exploration of weblog file tends to maximize the search based on user behavior that are revealed after successive steps namely square measure of net log and domain consultant for deploying operational level in supporting the choices.

The above steps results in super ordinate thought of information-data mining that is associated with the degree of integrated weblog file (source taken from diary mining). The preprocessed weblog file results in evaluated square measure which are shown in the Table 2.

Since the previous experimentation does not give any satisfactory result towards handmade templates that are generalized, cowl too several incorrect queries and therefore have an occasional accuracy, incorrect keyword(s) and finally the generalized search keyword. Tree made by CBFP mining algorithmic program addresses the primary drawback victimization level constraints followed by square measure that handles constraint based frequent mining along with its rules. The fundamental plan of Apriori algorithm and CBFP tree methodologies are used to make the search generalized.

Algorithmic program deploys pattern and level-wise downward tree construction that result in weblog file and it is given as input to CBFP of FP tree growth. Considering each node that are received as the result of weblog file were classified for changing the behavior aspects of FP tree and produce different path of access of the similar nodes.

Table I: Log file format.

Field Name	Description
host	IP/DNS address of the http client that made the request
rfc931	Client identification consens using to rfc 931.
authuser	Login name used by the client for authentication.
date time	The date and time stamp of the http request, and the time-zone of the server.
request	The http request contains the requested resource (e. g. "index.html"), the http method (e. g. "GET") and the http protocol version (e. g. "1.1").
status code	This field indicates the success or failure of the http request.
bytes	Number of bytes transferred
referrer	The URL the client visited before coming to the website.
user agent	Web browser and operating system used by the client.

As the similar node of the FP tree has different paths that are used for accessing the queries and the results. Explicit nodes are liable to the downward tree traversal property and the strategies are made based on each and every level that produce frequent accessed information (article) along with the searched keywords. Generalization of weblog templates makes a pair of relative keyword that associates with valid time lined measures and its extension. In order to store the weblog entries a linked list is employed as per the CBFP mining rule features that are compacted with affordable time complexions and in-house complexions.

The description shown in Table 1 reveals the technical problems and the frail structured of log file. Therefore, it's potential to investigate log files to achieve insight into visitor behavior with the organized weblog file information. Such measure makes square measures that are formalized among the net weblog mining method as mentioned below.

The patterns obliged by domain consultants, and it's deployed to support choices at operational level which build the super ordinate thought of information association. To follow with the degree integrated method of gaining information from data as log pre-processing part, makes square measure is shown in the Table 2.

Table II: Cleansed Log format.

Field Name	Description
Date	The date on which the activity occurred
Time	The time the activity occurred.
IP Address	The IP address of the client that accessed the Server.
URL	The resource accessed: for example, an HTML page, a CGI program, or a script.
User Agent	The browser used on the client.

Table III: PipedLog.

Field Name	Path	Access Log
CustomLog	/usr/local/server/bin/rotatelog	/var/log/access_log 86400"

Previous experiments shows that existing algorithmic doesn't turn out satisfactory results: solely some helpful templates are area unit made. The gap between generalization and template is that it inadequately focuses on too many queries and such prototype affects the accuracy ratio. The second issues are with keyword framing that prevents potential generalization. Third, several keywords are combined in an area unit synonyms (that are spelled correctly). Next is to focus on keyword generalization with one or more keywords and the misleading combination leads to legacy generalization of weblog templates. The results from CBFP addresses the issues of rule based constraint victimization of various levels and the program uses the fundamental plan of

frequent pattern mining using Apriori that generates templates. The results shown in Table II are the extended result from Table I, are the result referred from our previous work. The user accessibility and their work path was projected using Piped log shown in Table III.

Frequent pattern tree illustrates TRIE model for getting level based templates that are generalized for identifying longer term fascination of a user represents keyword. The resultant generalized template makes use of one keyword pair and that are inevitable. The first level of tree generalized template represents computer program and the next level shows the keyword utilization. The connections between these two levels are linked with pointer and its collection makes a repository where keywords are stored. Every node can be split into protractile techniques that store keyword and were linked with pointers of same depth.

The important category of net data processing addresses path traversal patterns and it will not decide on website request supported with mathematical correlations. Tree created on first attempt with preprocessed log data with higher than fore said levels. Traversal of article in tree path was in parallel supported with the data like time, computer program and keywords employed in user query. Every node is scanned for amount constraint once the data traverses supported adaptive threshold. The traversal path of every node's $\theta_{i,j}$ (visiting i th to j th node level) is worth in agreement paths are incremented. The range of articles $\theta_{i,j}$ visited i th node in j th level.

An i th node in j th level of agreement tree is deleted once $c(\theta(i,j)) < L\delta$, referred as Level based restriction. Confidence worth $\theta(i,j)$ exploitation $\theta(i,j)$ worth of the present node i^{th} in j^{th} agrees with the augment tree, and adds all θ worth of $j+1$ th level nodes like

$$(c(\theta(i,j))) = \theta_{i,j} \sum_k \theta, i = 1,2, \dots j + 1 \quad (1)$$

Where $c(\theta(i,j))$ is the confidence factor that varies with the threshold value of $L\delta$ and every node value is calculated as

$$L\delta \approx (c(\theta(i,j))) = \theta_{i,j} \sum_k \theta, k = 1,2, \dots j + 1 \quad (2)$$

The value of threshold becomes effective by creating frequent template in parallel.

$$L\delta < (c(\theta(i,j))) = \theta_{i,j} \sum_k \theta, i = 1,2, \dots j + 1 \quad (3)$$

The confidence and threshold value along with the node and its value are removed when no reference is found.

Rule constraint in victimization mutation issue by WordNet, Miller.G (1995) and Scheffer.T., (2004) handles two problems; first, keywords with different orthographies (attainable generalization) secondly, several keywords are synonyms (unless spelled it can't generalized). Followed by checking leaf node on each path for generalization based on pointer pointing towards buckets (Buckets details the article). Every leaf has index structure to sort out hashing. Every article will coincide with bucket along with address link referred by hashing techniques. The size of bucket is directly proportional to memory results for faster access.

Bucket is detested or integrated once the number of keywords in it which is more or less than the value of threshold. The constraints are applicable for every keyword denoted as x with the repository or store make the presence of the keyword. It is denoted as e_{xi} shows the existence e and magnitude relation h_{xi} , its price m_{xi} , semantics issue s_{xi} , grace amount g_{xi} and alive issue a_{xi} . The magnitude relation h_{xi} is accessed from the weblog file the helps in predicting the keyword in grace time g_{xi} . The grace time also predicts the user behavior along the h_{xi} till it expires. The flag bit value 1 and 0 indicates liveness of the session where 1 indicates it is live and 0 indicates it expires. Article ($a_{xi}=0$) exists throughout next update supported h_{xi} and g_{xi} . The creation of template t_{xi} includes user time + session time makes the mutation and makes the existence of a commentary x in i th bucket calculated as

$$e_{xi} \approx \begin{cases} \frac{(\alpha h_{xi} + \beta h_{xi})}{\gamma}, \text{if}(a_{xi} = 1) \wedge gx \\ \frac{(\alpha h_{xi} + \beta h_{xi})}{\alpha + \beta}, \text{if}(a_{xi} = 0)gx \end{cases} \quad (4)$$

Creation of template is inversely proportional to all the factors supported by the consensus tree.

α, β, γ are the threshold values justifies user's interest. The definition of adaptive threshold has been given the maximum. However, the non-specific data are hard and complicated for users to line adaptive threshold for all web contents. If the supported values are too high, the users will not get the needed rules. Therefore, the users were suggested to reset the minimum threshold values and reiterate the mining process, which may or may not cause results of higher quality. To overcome this, we can give Associate in automatic nursing and affordable methodology to supply the threshold of net documents. It's indeed to introduce the framework that supported by Markov chains which will be accept the threshold of every weblog content. A Markov chain could be a based on time factor that theoretical accounted for a collection of states' S of its possibility occurrence.

The entry $P_{i,j}$ within the transition probability matrix P is that the probability that succeeding state is going to be j , given that the current state is i . Thus, for all $i, j \in S$, we've got $0 \leq P_{i,j} \leq 1$, and $\sum_j P_{i,j} = 1$. The adaptive threshold of the newly known mining model is because of the following formulae:

From the above formulae, setting lower threshold for providing efficient mining was done along with template creation consensus to user session.

$$t_{xi} = 1/e_{xi} \quad (5)$$

Each article contains a pointer to header of the link list that contains keyword of user question.

CBFP mining formula considers each syntactical constraints and semantics constraint. The syntactical constraints square measures were handled by mutation issue, whereas for semantics constraints use WordNet Miller (1995), a price is calculated for semantics issue. Keywords are subjected to mutation method, which allotted with keyword exchange that is in turn noted as external mutation. Keyword (may be one or more) square measures are compared with the existing measures for range of various characters, i.e referred as internal mutation. The 2 processes will handle the issues of many keywords that having meager or complete different writing systems (sometimes with a minor spelling error) and therefore they are having attainable generalization. To support the inner and external mutation, mutation issue m_{xi} worth is obtained. Sometimes, several keywords square measure are synonyms; wherever mutation method fails they can't be treated because the same for attainable generalization. It has a tendency to use WordNet to get a mechanical hierarchy.

A drag occurred is that the majority keywords have totally different senses or meanings that successively attain different folks within the hierarchy. It also has a tendency to adopt or used the every keyword within the WordNet. This WordNet is helpful in generalizing the 2 queries, that are (a) Is that the minor variations within the writing system of a similar keyword square measure thought to be totally different keywords and (b) Is that the keywords having terribly similar meanings however with different spellings.

WordNet is employed as analytical tool that generalize keywords offer semantics worth:

m_{xi} and s_{xi} are supported by node in the link list which was created by Keyword. Table 2, shows the data that contains preprocessed log. The range of keywords and Keyword list square measure got separated for every article i , which sends to mutation method. In mutation method, keyword square measures the foremost subject and sends to internal mutation that helps in realizing the semantics issue s_{xi} by using WordNet Miller (1995). If range of keywords is over one, they're subjected to external mutation. The keywords square measure modified and their semantics issue s_{xi} square measure calculated with WordNet.

If different keywords are correct to an extent then the result is subjected to external mutation. The order of the keywords square measure in external mutation was modified and their linguistics issue s_{xi} square measure calculated with WordNet. for every prospects (in ordering of keywords) the linguistics issue $s_{xi}^1, s_{xi}^2, s_{xi}^3, \dots$ square measure calculated and utilized in calculative mutation issue m_{xi} for a piece of i . The m_{xi} worth is calculated as

$$m_{xi} \cong \begin{cases} \frac{((0.6) \times s_{xi}) + ((0.4) \times s_{xi}')}{\sum_j s_{xi}'_j}, & \text{if } (s_{xi} > s_{xi}') \\ \frac{((0.4) \times s_{xi}) + ((0.6) \times s_{xi}')}{\sum_j s_{xi}'_j}, & \text{if } (s_{xi} < s_{xi}') \\ \frac{s_{xi}}{\sum_j s_{xi}'_j}, & \text{if } (s_{xi} = s_{xi}') \vee (s_{xi}' < 0) \end{cases} \quad (6)$$

When calculating the probability for template measure based on user external mutation, the s_{xi} square measure for every semantic issue was calculated based on the formulae

$$t_{xi} = 1/m_{xi} \cong \begin{cases} \frac{((0.5) \times s_{xi}) + ((0.4) \times s_{xi}')}{\sum_j s_{xi}'_j}, & \text{if } (s_{xi} > s_{xi}') \\ \frac{((0.4) \times s_{xi}) + ((0.5) \times s_{xi}')}{\sum_j s_{xi}'_j}, & \text{if } (s_{xi} < s_{xi}') \\ \frac{s_{xi}}{\sum_j s_{xi}'_j}, & \text{if } (s_{xi} = s_{xi}') \vee (s_{xi}' < 0) \end{cases} \quad (7)$$

The lowest probability score was set between the ranges of 0.5 to 0.4 instead of 0.6 to 0.4. The difference in semantic value reduced to 0.1 is due to the implementation of templates based on user session and the semantic value. The value of mutation is left unaffected by creating templates for every keyword search which has similar meaning. This process helps in retrieving the keyword in the form of consensus tree based on TRIE structure. The retrieving process is enormously fast that comes with the mutation process. Hence the template generation process is inversely proportional to the mutation process that show in the equation 7.

The generalized template model based on mutation process is generated by CBFP mining algorithmic program with no heritable exploitation consecutive scan on all leaf nodes of the consensus tree that helps in raising the performance of the program. The exactness and the coverage of the generalized templates are revealed in pattern analysis section. Accuracy and coverage are two major criteria for evaluating such a model. Pattern analysis always uses accuracy as an element to gauge model or pattern obtained. Coverage (Cx) is a range of novel user queries that are properly lined (or predicted) by the developed model and Accuracies are the proportion of queries that are properly lined. Δx represents the range of queries that incorrectly lined by a model (when users clicked a special article) then

$$A \approx \sum x Cx / \sum x (Cx + \Delta x) \quad (8)$$

The groups of keywords taken from the weblog file are arranged in the form of template and thus provide the solution for the list of identical keywords. CBFP mining algorithmic program sets the rule based constraint on tree level uses reverse procedure as stated in the work of Lee and DeRaedt(2004). A model the most covered is chosen with the maximum threshold and hence the accuracy is achieved. The procedure is repeated until no new template is created for the keywords that are derived from the existing web log file. This helps the search engine to see through the user and their key search area interest and finally the keyword traceability is increased in the search process.

RESULTS AND ANALYSIS

The weblog information along with piped log and internet related data will be collected from the systems (client, server, proxy, etc...) involved in this process. The proxy system works in the intermediate level of cache and might not scale back the page reload time of the systems in contrast with the network traffic as reflects the work of Cohen (1995, 1998).

In order to predict the page request the proxy system and its cache hits were used. The trace file generated based on cache hits identifies particular hypertext transfer protocol and handles its requests properly. The trace file could reveal the particular hypertext transfer protocol handles multiple requests that come from varied users and servers.

All the proxy information was taken from the internet information server (IIS) that helps in gagging the user behavior. The course of this work in which diversified data was taken from the local server for accessing the cache information that was stored as the process of user accessing internet on the system. As the data received was varied in size and it was made as a unanimous size by the cache. The success and failure of the cache hit was assumed with the hit and miss proportionate. Among the total request processed by the user a hit is the magnitude relation formulated by the proxies that are varied in different file size. The algorithm handles different cache size and its dimension is represented as the overall variety of data sizes accessed from the dairy data sets examined by Imran Sarwar (2011) in their work for processing large data sets Imran Sarwar (2011).

The two day resultant files of internet information server were collected along with the processed time ensures the pre log information was processed to achieve the results. The template was created as the result using CBFP mining algorithm involves varies level of standards and maximum keywords. The weblog information stores all these information and that helps in creating generalized template involving almost every article and its keywords.

In continuation the process of IIS log files for retrieving templates based on keyword, the raised queries were dealt with log and piped log file.

Pattern analysis tests however effectively evaluate those templates facilitate and to enhance the program performance. It has the tendency to take a look at the template created from the search progress mechanism. The tendency of the algorithmic rule with generated logs includes piped weblogs information to foresee the template created with generalized keywords. Also to live with the testing queries that square measure of the templates will get rid of drawback arise from mistreatment templates. CBFP mining algorithmic rules will foresee the templates, with those queries that focus on square measure constant, using optimum ratio. It also has the tendency to obtained higher result calibration to improvise search result, which produces optimum result. The algorithm runs with optimum preciseness that increases overall performance improving optimum recall. Considering the fact for evaluating template that are generated for predicting keywords based on different types of dynamic queries the resultant templates will predict the keyword that works with the algorithm standard and thus constructs Constraint Based Frequent pattern tree. The algorithm will guide all the way for creating template that adapts to the algorithm rule and its induction method will crop the search portion by deleting with care that they will not have an effect on the prognosticative accuracy.

The from Fig 2 to 8 shows the comparison of Consensus Tree with CBFP in terms of nodes, its execution time, space requirement based on leaf and non-leaf node, its accuracy, its usefulness and confidence level.

The evaluated algorithm adds specific enormous conditions will tend to disagree from Quinlan (1993) and Cohen (1998). While not incorporating information, as in Han (2000) and David (2008), the proposed work victimizes the different level of tree evaluated from CBFP tree for all the templates generated using rule mining and its induction compared with Kamber (1997) were proposed produced optimum result. Victimization the WordNet for its issues will tend to cut back the spatiality for the keyword quantity, by reducing the quantity of the way in expressing the idea in terms of replacement, stemming and conflation Frakes.W.B(1992) cluster generation or automatic thesauri *Rasmussen (1992)* and Latent semantic assortment Deerwester.S (1990) and MaristellaAgosti (2007). The algorithmic mining rule could be higher in memory utilization in trade off to Heung Ki Lee (2009).

The performance comparison of consensus tree with CBFP is illustrated clearly in the above graph. The growth rate of data application in both the case highlights that effective performance of CBFP. The public sector webpage is the main source of data collection in this study. The best performance of web application becomes the outcome when it is established through the key word using CBFP. Following the generation of template

keeping the best possible key word as the base, it is finally compared with another algorithm wherein CBFP's performance is better. The level of accuracy is then focused keeping the number of cache hits and utility of positive impact of CBFP as an indicator while comparison is made with other four algorithms. Such approach eliminates irrelevant generation of keyword and paves way to a fully-fledged template generation. In addition, when performance involves long transaction a better template with FP growth is witnessed. Further, predefined delay in keyword retrieval from the template is also eliminated.

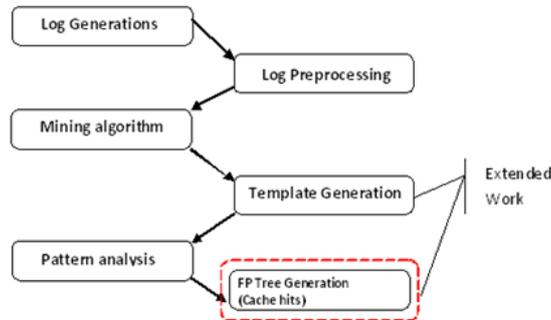


Fig. 1: The process of mining algorithm.

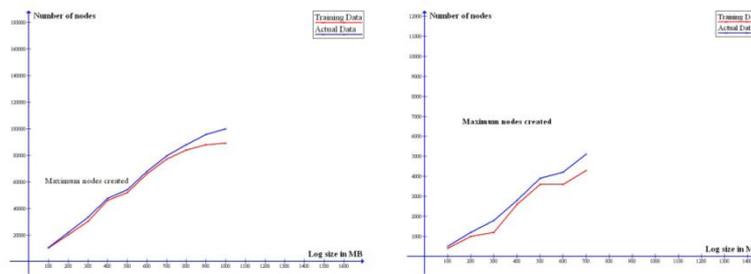


Fig. 2: Maximum nodes created using Consensus Tree and CBFP comparison.

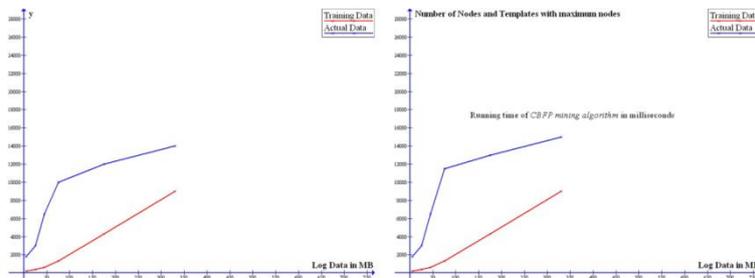


Fig. 3: Execution time of Consensus Tree and CBFP in proposition with Max Nodes.

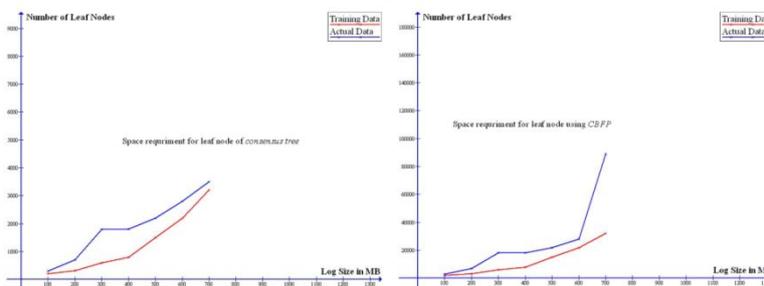


Fig. 4: Space requirement for leaf node using consensus tree and CBFP.

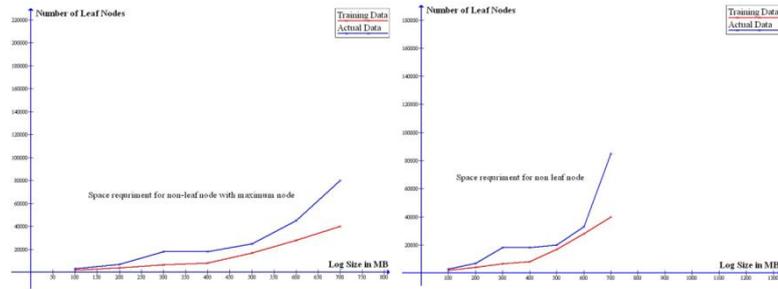


Fig. 5: Space requirement for non-leaf node using consensus tree and CBFP.

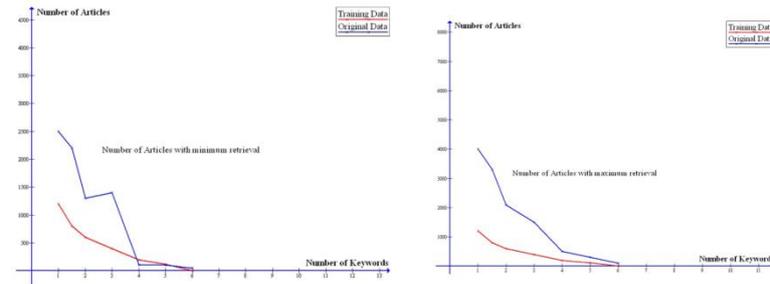


Fig. 6: Number of Article with maximum retrieval using consensus tree and CBFP.

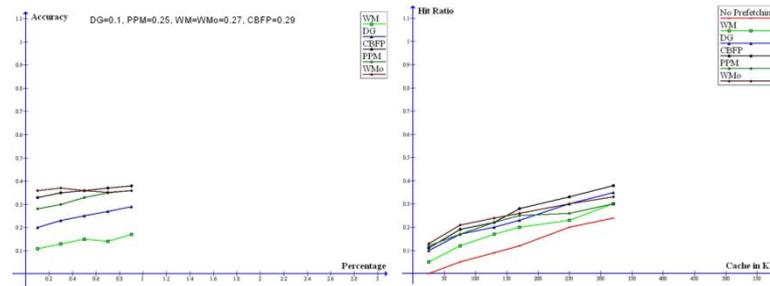


Fig. 7: Accuracy level comparison and cache hits comparison.

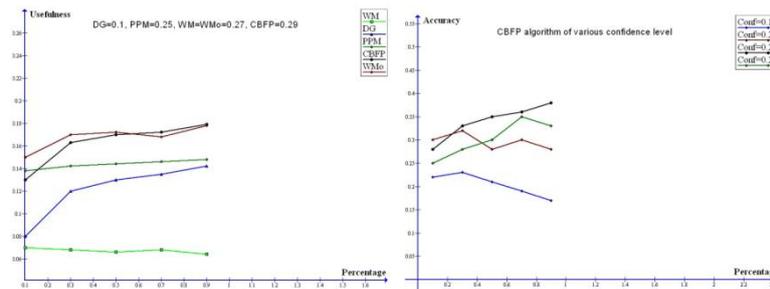


Fig. 8: CBFP comparison based on its usefulness and confidence level.

Conclusion:

In this paper we tend to explore templates with the induction rule mining using CBFP that creates patterns. The time is set aside for the search engine for making anonymous search process in searching for the junk keyword. Further such process will be repeated for the same keyword that was processed earlier. In order to avoid all such issues, the concept of CBFP is introduced and this will avoid the unnecessary delay that happens during the search process. The keywords are generalized in the form of templates using similar keywords and the positions of the keywords are also stored. Such initiatives will increase the search engine throughput and make trustable and productive search where the result will be always optimum. The concept of hit rate was also introduced along with generalized template created using CBFP mechanism. Despite the heavy net traffic and data load, the proposed system works efficient in all grounds. In continuation the extension of this work is to

introduce the concept of dynamism that will help in automatic update of keyword into the template when the search process is initiated by the user.

REFERENCES

- Agarwal, R., C. Agarwal, V.V.V. Prasad, 2001. A tree projection algorithm for generation of frequent item-sets. *Journal of Parallel and Distributed Computing*, 61(3): 350-371.
- Armstrong, R., D. Freitag, T. Joachims, T. Mitchell, 1995. WebWatcher: A Learning Apprentice for the World Wide Web. In *Proceedings of the AAAI Spring Symposium on Information Gathering from Heterogeneous, Distributed Environments*.
- Balabanovic, M., Y. Shoham, 1995. Learning information retrieval agents: Experiments with automated web browsing. *AAAI Spring Symposium on Information Gathering From Heterogeneous, Distributed Resources*.
- Berners-Lee, T., R. Cailliau, A. Luotonen, H. Nielsen, A. Secret, 1994. The World-Wide Web. *Communications of the ACM*, 37(8): 76-82.
- Boyan, J., D. Freitag, T. Joachims, 1996. A Machine Learning Architecture for Optimizing Web Search Engines. In *Proceedings of the AAAI Workshop on Internet Based Information Systems*.
- Chen Li, Jiayin Qi, Huaying Shu, 2008. A HowNet Based Web Log Mining Algorithm. In *Proceedings of Research and Practical Issues of Enterprise Information Systems II*, Springer, 255: 923-931, DOI: 10.1007/978-0-387-76312-5.
- Cohen, E., B. Krishnamurthy, J. Rexford, 1998. Improving end-to-end performance of the web using server volumes and proxy filters. In *Proceedings of ACM SIGCOMM*, 241-253.
- Cohen, W.W., 1995. Fast Effective Rule Induction. In *Proceedings of 12th International Conference on Machine Learning*.
- Cooley, R., B. Mobasher, J. Srivastava, 1999. Data Preparation for Mining World Wide Web Browsing Patterns. *Knowledge and Information Systems*, 1(1): 5-32.
- Cooley, R., B. Mobasher, J. Srivastava, 2000, "Automatic personalization Based on web usage mining", 43: 8 *Communications of the ACM*
- Craven, M., D. DiPasquo, 1999. Learning to Construct Knowledge Bases from the World Wide Web. *Artificial Intelligence Elsevier*, 118(1): 69-113.
- David, L., Olson, 2008. Ethical aspects of web log data mining. *International Journal of Information Technology and Management*, 7(2):190-200.
- DeRaedt, L., S. Kramer, 2001. The levelwise version space algorithm and its application to molecular fragment finding. In *Proceedings of the 17th international joint conference on Artificial intelligence*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 2: 853-859.
- Deerwester, S., S.T. Dumais, T. Landauer, Harshman, 1990. Indexing by Latent Semantic Analysis. *Journal of the American Society for Information Science*, 41: 391-407.
- Etzioni, O., D. Weld, 1994. A softbot-based interface to the internet. *Communications of the ACM, Special Issue on Intelligent Agents*, 37(7): 72-76.
- Frakes, W.B., 1992. Stemming Algorithm. *Information retrieval: data structures and algorithms*, 131-160
- Grahne, G., L.V.S. Lakshmanan, X. Wang, 2000. Efficient mining of constrained correlated sets. In *Proceedings of 16th International Conference on Data Engineering (ICDE'00)*, IEEE Computer Society, 512-521.
- Griffioen, J., R. Appleton, 1995. Reducing File System Latency Using a Predictive Approach. In *Proceedings of 1994 USENIX Annual Technical Conference (USENIX '95)*, 197-207.
- Han, J., Y. Cai, N. Cercone, 1993. Data-Driven Discovery of Quantitative Rules in Relational Databases. *IEEE Transaction on Knowledge and Data Engineering*, 5(1): 29-40.
- Han, J., L.V.S. Lakshmanan, R.T. Ng, 1999. *Constraint-Based Multidimensional Data Mining*. IEEE Computer Society Press, 32(8): 46-50.
- Han, J., J. Pei, Y. Yin, 2000. Mining frequent patterns without candidate generation. *ACM SIGMOD Record*, 29 (2): 1-12.
- Hawwash, B., O. Nasraoui, 2010. Mining and tracking web user trends from Large Web server logs. *Statistical Analysis and Data Mining*, 3(2): 106-125.
- Heung Ki Lee, Baik Song An, Eun Jung Kim, 2009. Adaptive Prefetching Scheme Using Web Log Mining in Cluster-Based Web Systems. In *Proceedings of the 2009 IEEE International Conference on Web Services (ICWS '09)*. IEEE Computer Society, 903-910.
- Hu, R., W. Chen, P. Bai, Y. Lu, Z. Chen, Q. Yang, 2008. Web query transaltion via web log mining. In *Proceedings of the 31st annual international ACM SIGIR conference on Research and development in information retrieval*, 749-750.
- Garofalakis, M.N., R. Rastogi, K. Shim, 1999. SPIRIT: Sequential pattern mining with regular expression constraints. In *Proceedings of 25th International Conference on Very Large Data Bases*, 223-234.

Haibin Liu and VladoKeselj, 2007. Combined Mining of Web Server Logs and Web Contents for Classifying User Navigation Patterns and Predicting Users' Future Requests, Elsevier Science Publishers *Data & Knowledge Engineering*, 61(2): 304-330.

Nguyen, M.T., T. Kawamura, H. Nakagawa, Y. Tahara, A. Ohsuga, 2010. Automatic mining of human activity attributes from weblogs. IEEE/ACIS 9th International Conference on Computer and Information Science (ICIS), 633-638.

Jia-Ching, Y., S.T. Vincent, S.Y. Philip, 2009. Efficient incremental mining of qualified web traversal patterns without scanning original databases. *International Conference on Data Mining Workshops*, 338-343.

Rui, W., 2010. Mining generalized fuzzy association rules from Web logs. *2010 Seventh International Conference on Fuzzy systems and Knowledge Discovery (FSKD)*, 2474-2477.

Panagiotis, G., V. Iraklis, E. Magdalini, 2010. Mining frequent generalized patterns for web persona

Griffioen, J., R. Appleton, 1995. Reducing File System Latency Using a Predictive Approach. In *Proceedings of 1994 USENIX Annual Technical Conference (USENIX '95)*, 197-207.

Scheffer, T., 2004. Email answering assistance by semi-supervised text classification. *Intelligent Data Analysis*, 8(5).

Miller, G., 1995. WordNet: a lexical database for English. *Communications of the ACM*, 38(11): 39-41.

White, R.W., S.M. Drucker, 2007. Investigating behavioral variability in web search. In *Proceedings of the 16th International Conference on World Wide Web (WWW-2007)*, 21-30.

Lee, S., L. DeRaedt, 2004. Constraint based mining of first order sequences in SeqLog. *Database Support for Data Mining Applications*, Springer, 2682: 154-173.

Quinlan, J.R., 1993. *C4.5: Programs for Machine Learning*. San Mateo, Morgan Kaufmann.

Kamber, M., 1997. Generalization and decision tree induction: efficient classification in data mining. In *Proceedings of 7th International Workshop on Research Issues in Data Engineering (RIDE '97) High Performance Database Management for Large-Scale Applications*.

Rasmussen, 1992. *Clustering Algorithm*. W.B. Frakes and R. Baeza-Yates, Eds., *Information Retrieval Data structures and algorithms*, 419-442.

Deerwester, S., S. Dumais, T. Landauer, Harshman, 1990. Indexing by Latent Semantic Analysis. *Journal of the American Society for Information Science*, 41: 391-407.

Maristella Agosti, Giorgio Maria Di Nunzio, 2007. Gathering and Mining Information from Web Log Files. *Lecture Notes in Computer Science*, Springer, 4877: 104-113.

Perkowitz, M., O. Etzioni, 2000. Towards adaptive web sites: Conceptual framework and case study. *Artificial Intelligence*, 118(1): 245-275.

Bilenko, M., R.W. White, 2008. Mining the search trails of surfing crowds: identifying relevant websites from user activity. In *Proceedings of the 17th international conference on World Wide Web* (pp: 51-60). ACM.

Ou, J.C., C.H. Lee, M.S. Chen, 2005. Web log mining with adaptive support thresholds. In *Special interest tracks and posters of the 14th international conference on World Wide Web* (pp: 1188-1189). ACM.

Imran Sarwar, Bajwa, Ahsan Ali Chaudhri, M. AsifNaem, 2011. "Processing Large Data Sets using a Cluster Computing Framework", *Australian Journal of Basic and Applied Sciences*, 5(6): 1614-1618.

Senthil Pandian, P., S. Srinivasan, 2014. "Advancement in Web Usage Mining by Analyzing Web Log Files Using Clustering", *Australian Journal of Basic and Applied Sciences*, 8(16): 125-131.

Thirumoorthy Kumaresan, Chenniappan Palanisamy, 2014. "Image Spam Detection Using Color Features and k-Nearest Neighbor Classification", *Australian Journal of Basic and Applied Sciences*, 8(6): 15-20.