

Pattern Discovery Techniques in Web Mining: A Review

Anirudh Chauhan¹, Neeraj Chauhan²

1 CCSIT, TMU

MORADABAD, INDIA

2 ASSOCIATE PROFESSOR

CCSIT, TMU

MORADABAD, INDIA

¹vishalrajput7107@gmail.com

Neeraj.computers@tmu.ac.in

Abstract— Web is the largest assortment of evidence; users might devote extra time period on the web outcome the apposite data or services they are troubled with. The data available is in form of structured (relational) and text data. Therefore, different kinds of data model can be implementable with web data for pattern discovery. Web mining is a data mining tool where the web related data is evaluated for pattern discovery and user navigation pattern. Additionally, according to the nature of data, the kind of mining is also upgraded. Pattern discovery is used to make a Web site additional approachable to the exclusive and specific desires or requirements of each individual user or set of users. This paper affords a detailed analysis of various approaches of pattern discover in data mining based on different domains with their compensations and limitations. A brief comparison has been made between the different techniques based on confident restrictions.

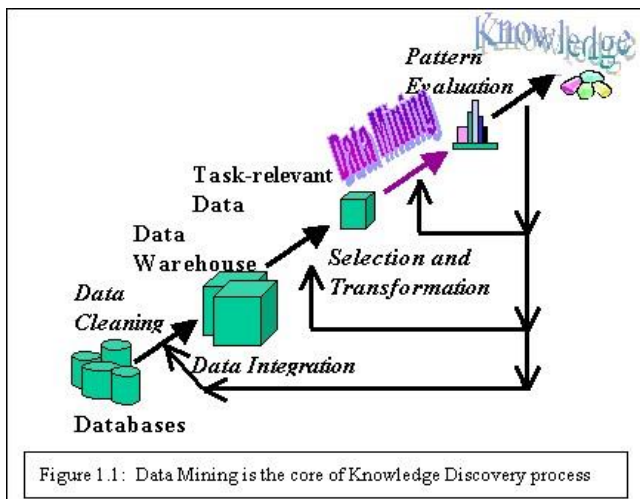
Keyword— web mining, web usage mining, log analysis, data models, pattern discovery.

I. INTRODUCTION

The internet is flooded with a lot of useful and useless information. It is very hard to define useful evidence for a particular user which is varying from time to

time. The useful information of one particular time may not be useful on different time or a different situation. The web itself is concerning day by day with newer technologies. Since internet is consuming free style medium that consents structured, non-structured, well-organized, nonordered format to provide an information in the web, finding not only the related information but to plan them according to the attention of a user is also a key challenge today and is known as web personalization.

“Pattern discovery is the approach of tailoring a website or content of website to the desires of every individual user or group of users or group, taking advantage of the information and web services attained through the investigation of the customer’s navigational concert”. Pattern Discovery is used to provide services to each specific user in a tailored manner. The broadcast of statistics on the internet has ended the personalization system an commitment. The discovered method must have the capability to resolve the supplementary data complications and let the customers run-through at least exertion to find the data they require.



When the customers work on the website their actions can be considered into two wide sets: browsing and searching. If the customers want to search data on internet, they want to deliver the web scheme with a search demand. If the customers don't provide the web system a precise search request, the system would reoccurrence huge volume of in appropriate data which can't fulfill the requirement of the user. If the customers' request is explicit, searching data on the web might turn out to be simple. To deliver better recommendation to the user the web personalization system's developers must identify what the customers' behavior when they surf on the web.

.Web usage mining:the web usage mining allows finding patterns from Web access information. This usage data provides the paths and user access patterns leading to accessed Web pages. This information is often gathered inevitably via the Web servers.

Web content mining: the web content mining is also known as text mining. In content mining applications the scanning and mining of text, pictures and graphs of a Web page is performed. That may help to conclude the consequence of the content.

Web structure mining: web structure mining is a tool, which is used to distinguish the connection between web pages. This organization of data is discoverable by the condition of web structure schema through database procedures for Web pages. This kind of data analysis allows a search engine to pull data regarding to a exploration query directly to the concerning Web page from the Web sites.

II. FORMATS OF DATA

The web access information can be establish in different places. Between origin servers to the client end, this access statistics is organized in different arrangements that are itemized in this section.

Proxy Servers: A proxy server is a software system. That is basically implemented by an organization that is connected to the Internet. Therefore, proxy servers are acts as an intermediary between a host and the Internet connectivity. Using this application the concerning organization can ensure security, reserving services and administrative control. Proxy servers can also be a respected source of usage data. A proxy server also manages access logs, in similar format to Web servers, this access log help to record Web page requests and retorts from the web servers.

Client Side Data: Client side data are serene from the host. That is currently accessing the Websites. To collect information directly from the client end, such as the time that the user is retrieving and leaving the Web site, a list of sites visited before and after the current site a client agent may accommodating. Client side data are more dependable than server side data. On the other hand, the use of client side data acquisition methods is also inspiring. The main delinquent is that, the different agents accumulating information.

Cookies: In addition to the use the web access log files, a different method often used in the collection of data is the tracking of cookies

Server Log Files: Server side data are together at the Web servers of a web site. The web server automatically generates the log file when a user request is made

from that server. These logs store Web pages information that is retrieved by the companies of the site. Most of the Web servers support as a default option the Common Log File Format, which includes information about the IP address of the client making the entreaty, the host name and user name, the time stamp of request, file name that is demanded, and the file size. The Extended Log Format (W3C), which is supported by Web servers such as Apache and Netscape, and the similar W3SVC format, supported by Microsoft Internet Information Server, include additional information such as the address of the mentioning URL to this page, i.e., the Web page that carried the visitor to the site, the name and version of the browser used and the operating system of host machine.

III. RELATED WORK

1. Category Based

In category based, there are two approaches, the first approach is collaborative filtering patterns, and this permits customers to take benefit of other customers' interactive actions based on a degree of likeness between them. Another approach is ruled based pattern; rather than matching customers' response to the web content or summaries of other customers, this model match that inquiry to some fixed rules or settlements, about customer performance.

The system logger is intended to wrinkle customers' net convention information. Log files gather customers' visiting counts on every hyperlink on the Web pages. By applying some data mining technique we can excerpt data from the log. Category Generator can categorize the customers into different groups on the basis of log data. Category Generator can identify which customer belongs to which group. In this approach, author proposed a technique of pattern discovery system stretched from the exploration of classical personalization schemes

and associated knowledge. A novel system logger is intended in this approach to store all the content or item openly retrieved by an individual customer, and with the help of category generator, which splits the content into various categories and provide the most appropriate result to the user.

2. Fuzzy Logic Based

Two information segments and two processing segments are considered in this approach. Information module collect the user and service information and processing modules are used to measuring client liking and product filtering, which are the key mechanisms in this personalization method.

The preference learning is reinforced by the fuzzy logic method which deals with the vague data or facts from the user's happenings. The suggested method provides another concept for personalization that adds fuzzy logic for measurement of users' likeness. Fuzzy sets are described as a numerical method to signify and deal with ambiguity or unsure in this area which is provisional on membership functions. The membership function pronounces how each separate fact from input mapped to a affiliation value in the interval. The anticipated method deals with the vagueness of users' activities; the suggested system produced the most appropriate and meaningful value based on the user's comportment and their access time. To produce appropriate membership functions for fuzzy logic is big stimulating issues in fuzzy systems strategy. It is problematic task because it openly affects the correctness of fuzzy logic method.

Algorithms

1. Apriori Algorithm

It searches for large item sets during its initial database pass and uses its result as the seed for discovering other large datasets during ensuing

passes. Rules having a support level above the least are called large or frequent item sets and those below are called small item sets. The algorithm is based on the large item set property which states: Any subset of a large item set is large and any subset of frequent item set must be frequent. The first algorithm for mining all frequent item sets and strong association rules was the AIS algorithm. Shortly after that, the algorithm was upgraded and renamed Apriori. Apriori algorithm is, the most classical and important algorithm for mining frequent item sets. The Apriori algorithm performs a breadthfirst search in the search space by causing candidate $k+1$ -itemsets from frequent k item sets. The frequency of an item set is computed by counting its occurrence in each transaction. Apriori is a significant algorithm for mining frequent item sets for Boolean association rules. Since the Algorithm uses prior acquaintance of frequent item set it has been given the name Apriori. It is an iterative level wise search Algorithm, where k item sets are used to explore $(k+1)$ item sets the finding of each L_k requires one full scan of database. There are two steps for sympathetic that how L_{k-1} is used to find L_k :- The join step:- To find L_k , a set of candidate k item sets is generated by joining L_{k-1} with itself. This set of candidates is denoted C_k . The prune step:- C_k is a superset of L_k , that is, its members may or may not be frequent, but all of the frequent k -itemsets are included in C_k . A scan of the database to determine the count of each candidate in C_k would result in the determination of L_k . C_k , however, can be huge, and so this could involve hefty computation to reduce the size of C_k .

2. FP-Tree

A frequent-pattern tree is a tree structure. It consists of one root labeled as "null", a set of itemprefix subtrees as the children of the root, and a frequent-item-header table. Each node in the itemprefix subtree consists of three fields: item-name, count, and node-link, where item-name registers which item

this node represents, count registers the number of transactions represented by the portion of the path reaching this node, and node-link links to the next node in the FP-tree carrying the same itemname, or null if there is none.

FP-growth is a well-known algorithm that uses the FP tree data structure to achieve a condensed representation of the database communications and employs a divide and-conquer approach to decompose the mining problem into a set of the above problem by Reducing passes, Shrinking number of candidates and facilitating support counting of candidates. An FP-tree-based patternfragment progress mining method is developed, which starts from a frequent length-1 pattern (as an initial suffix pattern), examines only its conditionalpattern base (a "subdatabase" which consists of the set of frequent items co-occurring with the suffix pattern), constructs its (conditional) FP-tree, and performs mining recursively with such a tree. The FPgrowth algorithm is one of the fastest approaches for frequent item set mining. The FP-growth algorithm uses the FP-tree data structure to achieve a condensed demonstration of the database transaction and employees a divide-and conquer approach to putrefy the mining problem into a set of smaller problems. In essence, it mines all the frequent item sets by recursively finding all frequent item sets in the conditional pattern base which is efficiently constructed with the help of a node link structure.

A prefix tree is a data structure that provides a

compact representation of transaction data set. Each node of the tree stores an item label and a count, with the count representing the number of transactions, which contain all the items in the path from the root node to the current node. The frequent items are computed as in the Apriori algorithm and represented in a table called header table. Each

record in the header table will contain the frequent item and a link to a node in the FP-Tree that has the same item name. Following this link from the header table, one can reach all nodes in the tree having the same item name. Each node in the FP-Tree, other than the root node, will contain the item name, support count, and a pointer to link to a node in the tree that has the same item name. **The main components of FP tree**

It consists of one root labelled as "root", a set of item prefix sub-trees as the children of the root, and a frequent-item header table.

- Each node in the item prefix sub-tree consists of three fields: item-name, count, and node-link, where item name registers which item this node represents, count registers the number of transactions represented by the portion of the path reaching this node, and node link links to the next node in the FP tree carrying the same item-name, or null if there is none.
- Each entry in the frequent-item header table consists of two fields, (1) item-name and (2) head of node link, which points to the first node in the FP-tree carrying the item-name.

IV. COMPARATIVE ANALYSIS

In this section, we have debated comparison between above approaches based on different parameters. We have focused on techniques in corresponding approach and advantage

Algorithm/parameter	Advantage	Limitation
Apriori Algorithm	Searches for large item set	Full scan require for single item set

FP-Tree	Scan frequently	More complex for non-frequent item
Category Based	User can take the benefit of other users similar interest	Rule based result depend upon developer perception
Fuzzy Logic Based	Deal with uncertainty and ambiguity for better result	Correctness of fuzzy system is difficult

Table.1 Comparative Analysis.

V. CONCLUSION

We have discussed algorithm and approaches for pattern discovery based on different domains. They have some strength and weaknesses, but the motive of these work are to make more accurate Web recommendation and provide relevant information and services to each individual user at different point of time by these systems. Some methods are based on content of the web page and users' interest and some of the algorithm includes clustering and data mining techniques. A comparative analysis on the basis of certain parameters a brief comparison is being provided among the all discussed approaches.

ACKNOWLEDGEMENT

I would like to express my special thanks of gratitude to our principal **DR.R.K.Dwivedi** and my teacher **Miss.Neeraj Chauhan** who gave me the golden opportunity to do this wonderful research paper on the topic **PATTERN DISCOVERY TECHNIQUES IN ONLINE DATA MINING**, which also helped me in doing a lot of research and I came to know about so many new things I am really thankful to them. Secondly I would also like to thank my parents

and friends who helped me a lot in finalizing this research paper within the limited time frame.

REFERENCES

- [1] N. Sael, A. Marzak, and H. Behja, —Web Usage Mining Data Preprocessing and Multi Level Analysis on Moodle, IEEE, 2013. [2] N. Lakshmi, R. S. Rao, and S. S. Reddy, —An Overview of Preprocessing on Web Log Data for Web Usage Analysis, International Journal of Innovativ Technology and Exploring Engineering (IJITEE), Volume-2, Issue-4, March 2013.
- [3] H. peng, —Discovery of Interesting Association Rules on Web Usage Mining, International Conference. 2010. [4] J. Han, J. Pei, and Y. Yin, —Mining frequent patterns without candidate generation: A frequent-pattern tree approach; Data Mining and Knowledge, 2003.
- [5] R. Agrawal and R. Srikant, "Fast Algorithms for Mining Association Rules", VLDB. Sep 12-15 1994, Chile, 487-99, pdf, ISBN 1-55860-153-8.
- [6] S. Kumar and K.V. Rukmani, —Implementation of Web Usage Mining Using APRIORI and FP Growth Algorithms, 2010.
- [7] J. Srivastava, R. Cooley, M. Deshpande, and P. Tan, —Web usage mining: discovery and applications of usage patterns from web data, Vol. 1, No.2, 2000, pp.12–23.
- [8] C.C. Lee and W. Xu, —Category-Based Web Personalization System, International Conference on Web Information Systems and Technologies, IEEE 2001, pp.1372-1377. [9] B.Hua, K. Wai Wong and C.C.Fung, —Fuzzy Logic Based Product Filtering for Web Personalization In E-Commerce, IEEE 2007. [10] S. Veeramalai, N. Jaisankar, and A.Kannan, —Efficient Web Log Mining Using Enhanced Apriori Algorithm with Hash Tree and Fuzzy in 2010.
- [11] J. Srivastava, R. Cooley, M. Deshpande, and P. Tan proposed —Web Usage Mining: Discovery and Applications of Usage Patterns from Web Data, 2000.
- [12] J. Han, J. Pei, R. Mao —Mining Frequent Patterns without Candidate Generation in 2004
- [13] Y. Li, C. Zhang, and J.R. Swan, —An Information Filtering Model on the Web and Its Application in Jobagent, Knowledge-Based Systems, vol. 13, no. 5, pp. 285-296, 2000. [14] Y. Li and N. Zhong, —Interpretations of Association Rules by Granular Computing, Proc. IEEE Third Int'l Conf. Data Mining (ICDM '03), pp. 593-596, 2003. [15] Y. Li and N. Zhong, —Mining Ontology for Automatically Acquiring Web User Information Needs, IEEE Trans. Knowledge and Data Eng., vol. 18, no. 4, pp. 554-568, Apr. 2006. [16] Y. Li, X. Zhou, P. Bruza, Y. Xu, and R.Y. Lau, —A Two-Stage Text Mining Model for Information Filtering, Proc. ACM 17th Conf. Information and Knowledge Management (CIKM '08), pp. 1023-1032, 2008.
- [17] H. Lodhi, C. Saunders, J. Shawe-Taylor, N. Cristianini, and C. Watkins, —Text Classification Using String Kernels, J. Machine Learning Research, vol. 2, pp. 419-444, 2002. [18] A. Maedche, Ontology Learning for the Semantic Web. Kluwer Academic, 2003.
- [19] C. Manning and H. Schütze, Foundations of Statistical Natural Language Processing. MIT Press, 1999.
- [20] I. Moulinier, G. Raskinis, and J. Ganascia, —Text Categorization: A Symbolic Approach, Proc. Fifth Ann. Symp. Document Analysis and Information Retrieval (SDAIR), pp. 87-99, 1996.
- [31] J.S. Park, M.S. Chen, and P.S. Yu, —An Effective Hash-Based Algorithm for Mining

Association Rules, Proc. ACM SIGMOD Int'l Conf. Management of Data (SIGMOD '95), pp. 175-186, 1995.

[22] J. Pei, J. Han, B. Mortazavi-Asl, H. Pinto, Q. Chen, U. Dayal, and M. Hsu, —Prefixspan: Mining Sequential Patterns Efficiently by PrefixProjected Pattern Growth, Proc. 17th Int'l Conf. Data Eng. (ICDE '01), pp. 215-224, 2001.

[23] M.F. Porter, —An Algorithm for Suffix Stripping, Program, vol. 14, no. 3, pp. 130-137, 1980.

[24] S. Robertson and I. Soboroff, —The Trec 2002 Filtering Track Report, TREC, 2002, trec.nist.gov/pubs/trec11/papers/OVER.FILTERING.ps.gz.

[25] S.E. Robertson, S. Walker, and M. Hancock-Beaulieu, —Experimentation as a Way of Life: Okapi at Trec, Information Processing and Management, vol. 36, no. 1, pp. 95-108, 2000.

[26] J. Rocchio, Relevance Feedback in Information Retrieval. chapter 14, Prentice-Hall, pp. 313-323, 1971.

[27] T. Rose, M. Stevenson, and M. Whitehead, —The Reuters Corpus Volume1—From Yesterday's News to Today's Language Resources, Proc. Third Int'l Conf. Language Resources and Evaluation, pp. 29-31, 2002.

[28] G. Salton and C. Buckley, —Term-Weighting Approaches in Automatic Text Retrieval, Information Processing and Management: An Int'l J., vol. 24, no. 5, pp. 513-523, 1988.

[29] M. Sassano, —Virtual Examples for Text Classification with Support Vector Machines, Proc. Conf. Empirical Methods in Natural Language Processing (EMNLP '03), pp. 208-215, 2003.

[30] S. Scott and S. Matwin, —Feature Engineering for Text Classification, Proc. 16th Int'l Conf. Machine Learning (ICML '99), pp. 379-388, 1999.

[31] F. Sebastiani, —Machine Learning in Automated Text Categorization, ACM

Computing Surveys, vol. 34, no. 1, pp. 1-47, 2002.

[32] M. Seno and G. Karypis, —Slpminer: An Algorithm for Finding Frequent Sequential Patterns Using Length-Decreasing Support Constraint, Proc. IEEE Second Int'l Conf. Data Mining (ICDM '02), pp. 418-425, 2002.

