



An Optimizing Method for Image Extraction with Partial Tree Alignment Algorithm

Gaganpreet Kaur *,

*Student of M. Tech,

Department of Computer Science Engineering,

Sri Guru Granth Sahib World University,

Fatehgarh Sahib, Punjab, India.

Usvir Kaur

Assistant Professor,

Department of CSE,

Sri Guru Granth Sahib World University,

Fatehgarh Sahib, Punjab, India.

Abstract: *With the explosion of the World Wide Web, a wealth of data on many different subjects has become available online. Usually, users retrieve Web data by browsing and keyword searching. But, these traditional methods have their limitations and disadvantages. Search engine helps to retrieve the relevant web sites based on the keyword specified by the user. It performs various operations such as crawling, indexing etc. It displays thousands of links as a result of the web search, but there are many road blocks that can make this process difficult or even impossible. So, the proposed system mainly aims to eradicate the disadvantages of search engines by exploring the contents of a web page to a maximum extent. It finds the exact keywords that match a page. When the search engine searches for web pages related to exact keyword, it can return only a few pages which are highly focused, specific and relevant to the topic. By this, the end-user gets the required information related to the search. We have concentrate our research on image extraction with Partial tree alignment algorithm and the Experiment shows that new approach is feasible and effective as compare to previous one which one is based on text Retrieval from the Search Engine.*

Keywords:

I. INTRODUCTION[1,2,3]

Web Mining is the extraction of interesting and potentially useful patterns and implicit information from artifacts or activity related to the World Wide Web. There are roughly three knowledge discovery domains that pertain to web mining: Web Content Mining, Web Structure Mining, and Web Usage Mining. Web Content Mining is the process of extracting knowledge from the content of documents or their descriptions. Web Structure Mining is the process of inferring knowledge from the World Wide Web organization and links between references and referents in the Web. Finally, Web Usage Mining, also known as Web Log Mining, is the process of extracting interesting patterns in Web access logs. Web Crawlers are programs that exploit the graph structure of the web move from page to page. The key motivation of the web crawlers has been to retrieve web pages and add them or their representations to a local repository. Such a repository may then serve particular application needs such as those of web search engine[1]. The problem of extracting data from a Web page that contains several structured data records. The Objective is to segment these data records, extract data items or fields from them and put the data in a database table. There are two algorithms for the data extraction i.e. Top-down, bottom-up algorithm. On the basis of these two algorithms, there is a development of Hybrid algorithm called Bi-Direction Data Extraction. It can be able to extract and discriminate the relevance of different repetitive information contents with respect to the user's visual perception of the web page[3].

II. WEB MINING [3]

Currently, the World Wide Web (or the Web for short) is a huge information source. Before the Web, finding information means asking other person or looking for it in some books or other kinds of text document. Now, if we need information about something, we can just open a web browser and search it in web search engine. The Web is also a popular communication media. People interact with each other via web forum or social network web site like Face book and Twitter. Finally, the Web is also an important channel for conducting business. Many companies have used the Web for product campaign or to open online store. Because of those important uses of the Web, many researches have been conducted to extract useful information from the Web. Web mining aims to discover useful information or knowledge from the web hyperlink structure, page content, and usage data.[3] Based on those primary kinds of data used in the mining process, web mining tasks can be categorized into three types: Web Structure Mining, Web Content Mining and Web Usage Mining.

III. DATA EXTRACTION[4]

Web Image Extraction systems are a broad class of software applications targeting at extracting information from Web sources like Web pages [2]. A Web Image Extraction system usually interacts with a Web source and extracts data stored in

it: for instance, if the source is a HTML Web page, the extracted information could consist of elements in the page as well as the full-text of the page itself. Eventually, extracted data might be post-processed, converted in the most convenient structured format and stored for further usage. Web Data Extraction systems and extensive use in a wide range of applications like the analysis of text documents at disposal of a company (like e-mails, support forum, technical and legal documentation, and soon), Business and Competitive Intelligence [4], crawling of Social Web platforms [5], Bio-Informatics [93] and so on. The importance of Web Data Extraction systems depends on the fact that, today, a large (and quickly growing) amount of information is continuously produced, shared and consumed online. Web Data Extraction systems allow to efficiently collect this information with a limited human effort. The availability and analysis of collected data is an infeasible requirement to understand complex social, scientific and economic phenomena which generated the information itself. So, for instance, collecting digital traces produced by human users in Social Web platforms like Face book, YouTube or Flickr is the key step to verify sociological theories on a large scale [6].

IV. Related Work

In the previous research The contents of the Web Page were extracted which includes the source code, hyperlinks, Meta tags and keywords. The weight age were assigned to words based on the tags and the ranking was also given by calculating the frequency. The links were displayed based on the keywords. The main goal of the Previous system is based on extracted keyword frequency; hyperlinks may be created in future to provide a convenient way for users to retrieve related documents[2]. In future, this method can be combined with search engine for optimizing it .It can be implemented in other languages. The Extraction of images may be included. Depth of searching can be extended for each linked page.

V. Purposed Work

Our Purposed approach is to extract the images from the web pages Instead of Retrieve the keywords of Text from the Search engine. Firstly enter the query and the relevant pages come on the top .The user may give the URL of the Web page to be tested as input. The information can retrieve from the Web pages. Tags, Words, keywords, Hyperlinks can be extracted. Hence, database table has been created which recorded the all terms. We have used same partial tree alignment algorithm ,but in recent research it was based on text retrieval. In our purposed work it is based on image extraction from web usage. We have followed a methodology as listed down and compare our work with the recent research in terms of frequency of keywords outcome.

VI. Methodology of Purposed Work

It describes the whole image extraction process in which we extract the various Images from Database.

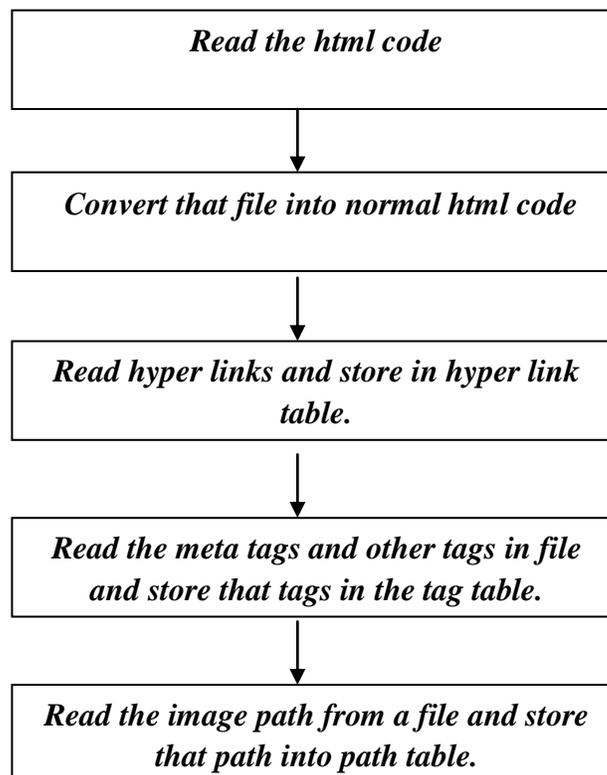


Figure: Methodology of Purposed Work

VII. Implementation Phase- Image Extraction Using Partial Alignment Algorithm

1. Weighted Page Rank
2. Segmentation
3. Tag Extraction
4. Content Extraction
5. Display content.

1. In first step, Weighted Page Rank algorithm (WPR) [6]: This algorithm is an extension of Page Rank algorithm. WPR takes into account the importance of both the in links and the out links of the pages and distributes rank scores based on the popularity of the pages. WPR performs better than the conventional Page Rank algorithm in terms of returning larger numbers of relevant pages to a given query. According to author the more popular web pages are the more linkages that other WebPages tend to have to them or are linked to by them. The proposed extended Page Rank algorithm—a Weighted Page Rank Algorithm—assigns larger rank values to more important (popular) pages instead of dividing the rank value of a page evenly among its out link pages. Each out link page gets a value proportional to its popularity (its number of in links and out links).

2. In the second phase, the Web page is split in segments, without extracting any data of images. This preprocessing phase is instrumental to the latter step. In fact, the system not only performs an analysis of the Web page document based on the DOM tree, but also relies on visual cues trying to identify gaps between data records. This step is useful also because helps the process of extracting structural information from the HTML document, in that situations when the HTML syntax is abused, for example by using tabular structure instead of CSS to arrange the graphical aspect of the page.

3. In the third step, the partial tree alignment algorithm is applied to data records earlier identified. Tag Extraction is the first module in the proposed system. It deals with extracting the tags automatically from the web pages. It requires identifying the HTML source of the web page then separating the tags and the content. Finally the tags are extracted separately. Each Tag is identified separately. The objective of this algorithm is to segment the data records, extract images items/fields from them and put the data in a database table. It consists of identifying individual data records in a page and aligning and extracting data items from the identified data records. Partial alignment aligns only those data.

4. Content extraction is the fourth module. It deals with extracting the contents from the Web pages. Content extraction is the main task. It is done along with the first module, Tag extraction. The Web page contains the information i.e. the data which is to be extracted, and these data are called as interesting data. The interesting information may also be called as the knowledge content of the Webpage. The content gives the details about the Web page. The content is built with various key words. Weight age is assigned to the tags. Each word is separated and the frequency of each word is calculated separately. Then the value of each word is calculated by using the formula $\text{Word Value} = (\text{frequency}) * (\text{weight age})$. The value calculation is based on the predefined weight age assigned to the tags. Then the priority is assigned to the key-words based on the highest value. This process is called as parsing. The noisy data present in the content such as and, where, when, though etc (stop words) are eliminated.

5. Display Content is the module deals with the user interface. It displays the page which is tested (given as input). It also displays the source code, content, links and the ranked Images related with them.

Algorithms Steps As: Let partial tree alignment (S)

1. Sort trees in S in descending order according to the number of Data items that are not aligned;
2. T_s = the first tree (which is the largest) and delete it from S;
3. Flag = false; R = 0; I = false;
4. While (S \neq 0)
5. T_i = select and delete next tree from S;
6. Simple_tree_matching(t_s , t_i);
7. L = align trees(t_s , t_i); // based on the result from line 6
8. If t_i is not completely aligned with t_s then
9. I = insert into seed(t_s , T_i);
10. if not all unaligned items in T_i are inserted into T_s then
11. Insert T_i into R;
12. end if;
13. end if;
14. if (L has new alignment) or (I is true) then
15. flag = true
16. end if;
17. if S = 0 and flag = true then
18. S = R; R = 0;
19. flag = false; I = false
20. end if;
21. end while;
22. Output data fields from each T_i to the data table based on the alignment results.

VIII. Result and Comparison

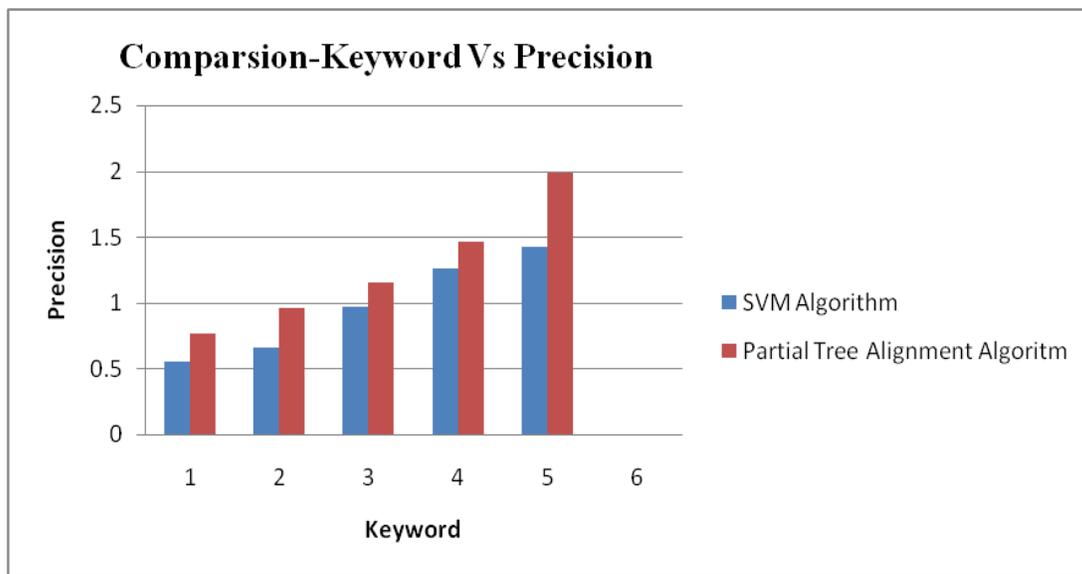
After the implementation of algorithm which is based on partial tree alignment we will extract images and their frequency of occurring in the data base with the help of search engine. In the Previous research we find that the search engine retrieve the frequency of words in text or images which was not true relevance of the retrieve records. So we have purposed an algorithm which was quiet efficient in terms of retrieval of images which are based along our search in the search engine. We have showed the comparison of both mechanism side by side and results are examined through graphs.

IX. Result Parameter

FREQUENCY OF WORDS: We have taken Frequency of images and words as a parameter of the comparison of our results and build up our comparison based on some sort of keywords listed in table given below:

Table1. Comparison of both Methods.

keyword		SVM Algorithm Series 1	Partial Tree Alignment Algorithm Series 2
1	PHONE	0.55	0.768
2	HOME	0.6565	0.96
3	PAGE	0.967	1.15
4	PRODUCT	1.26	1.46
5	ARTICLE	1.42	1.987



Graph Represents significant comparison Between Keyword and Precision.

X. Conclusion

The contents of the Web Page were extracted which includes the source code, hyperlinks, Meta tags and keywords. The links were displayed based on the keywords. The main goal of the proposed system is based on extracted keywords, Meta tags; hyperlinks may be created in future to provide a convenient way for users to retrieve related images. This method can be combined with search engine for optimizing it. Also this approach enables very accurate alignment of multiple data records. During this process no data items are involved, because partial tree alignment works only on tree tags matching, represented as the minimum cost, in terms of operations (i.e., node removal, node insertion, node replacement), to transform one node into another one. In the structure of the Web page at the time of the definition of the alignment. This implies that the method is very sensitive even to small changes that might compromise the functioning of the algorithm and the correct extraction of information. Even in this approach, the problem of the maintenance arises with outstanding importance.

XI. Future Work

This method can be combined with search engines for the optimizing it. It can be implemented in other languages ,so it can efficiently used by search engines to retrieve specified records.

References

- [1] Navjot Kaur, Jaspreet Kaur Sahiwal, Navneet Kaur, "EFFICIENT K-MEANS LUSTERING ALGORITHM USING RANKING METHOD IN DATA MINING", *International Journal of Advanced Research in Computer Engineering & Technology*, Volume 1, Issue 3, May 2012.
- [2] R. Baumgartner, W. Gatterbauer, and G. Gottlob. Web data extraction system. *Encyclopedia of Database Systems*, pages 3465-3471, 2009.
- [3] U. Irmak and T. Suel, "Interactive wrapper generation with minimal user effort," In Proc. 15th International Conference on World Wide Web, pages 553-563, Edinburgh, Scotland, 2006. ACM.
- [4] R. Baumgartner, O. Friolich, G. Gottlob, P. Harz, M. Herzog, P. Lehmann, and T. Wien, "Web data extraction for business intelligence the lixto approach," In Proc. 12th Conference on Datenbanksysteme in Biuro, Technik und Wissenschaft, pages 48-65, 2005.
- [5] S. Catanese, P. De Meo, E. Ferrara, G. Fiumara, and A. Provetti, "Crawling facebook for social network analysis purposes," In Proc. International Conference on Web Intelligence, Mining and Semantics, page 52, Sogndal, Norway, 2011. ACM.
- [6] L. Backstrom, P. Boldi, M. Rosa, J. Ugander, and S. Vigna, "Four degrees of separation," Arxiv preprint arXiv:1111.4570, 2011.
- [7] Wenpu Xing and Ali Ghorbani, "Weighted Page Rank Algorithm", In proceedings of the 2rd Annual Conference on Communication Networks & Services Research, PP. 305-314, 2004.
- [8] Tak-Lam Wong, Wai Lam, "An unsupervised method for joint information extraction and feature mining across different web sites", *Data & Knowledge Engineering*, Volume 68, Issue 1, January 2009, Pages 76-79.
- [9] Xiangwen Ji, Jianping Zeng, Shiyong Zhang, Chengrong Wu, "Tag tree template for Web information and schema extraction", *Expert systems with Applications*, Volume 37, Issue 12, December 2010, Pages 8492-97.
- [10] Gilles Nachouki, Mohamed Quafafou, " MashUp web data sources and services based on semantic queries", *Information Systems*, Volume 36, Issue 2, April 2011, Pages 151-173.
- [11] Viktor de Boer, Maarten van Someren, Bob J. Wielinga, "A redundancy-based method for the extraction of relation instances from the Web", *International Journal of Human-Computer studies*, Volume 65, Issue 9, September 2007, Pages 816-831.
- [12] Jer Lang Hong, Eu-Gen Siew, Simon Egerton, "Information extraction for search engines using fast heuristic techniques", *Data & Knowledge Engineering*, Volume 69, Issue 2, February 2010, Pages 169- 196.
- [13] Lirong Wan, Xinjun Wang, Congcong Chen, "A Sprial-Decoding Method for Web Data Extraction", *IEEE Conference Proceedings on First International conference on Education Technology and Computer Science*, 2009, Volume 1, Pages 1026-1029.