

Web Content Mining Tools for Information Extraction in Web Environment

Deven M. Kene

Dr. Pradeep K. Butey

Abstract - There has been a tremendous growth of the World Wide Web in recent years with millions of pages being added on a regular basis. Increasing amount of information are now using websites for multitude of tasks: Information search, Information extraction, Knowledge representation etc. However the flip side of the massive numbers of websites and content available online is that users often feel overwhelmed and intimidated by this information overload. In this scenario, Web content mining techniques is playing an important role in making the web experience easy and useful for all kinds of users. Web content mining Techniques involves many data mining techniques and help with both usability as well as user retention. This paper attempts to describe the various technique involved in web content mining for information extraction in web environment, how the web content mining process works and its benefits to users.

Keywords - Web Mining, Web Content mining, Web content Mining Techniques.

I. INTRODUCTION

Today's web has brought the revolution by transforming the developed world towards a knowledge economy. Its domain lies in vast areas such as information retrieval, extraction, integration, knowledge discovery, encompassing databases, artificial intelligence, etc. With exceedingly growth of the Web, there is an increasing volume of data and information published in numerous Web pages yielding some amount of noise in web pages representing information. Generally a web contains many kinds of information represented in the form of text, image, audio, video, etc. [1] and web pages in the form of headers, content, advertisements, panels of

navigation, copyright notices, footers etc. This form of information is application specific where in part of the information is useful for that application and other information become noise for other application. In web content mining some techniques have been used to extract and mine useful knowledge or information from these web pages.

Web mining is divided into three categories: Web content mining (WCM), web usage mining (WUM) and web structural mining (WSM). Web content mining is the process of identifying user specific data from text, image, audio or video data already available on web. Web pages may be structured, semi-structured or unstructured which necessitates the need of web content mining. The web content mining are used in information retrieval, knowledge extraction, resource finding, generalization, personalization, mining unstructured and semi structured content of web pages [2] . The Web content mining targets the Knowledge discovery, in which the main objects are the traditional collections of text documents

and the collections of multimedia documents such as images, videos, audios which are embedded in or linked to the web pages. Techniques belonging to the Web Content Mining such as classification and clustering, separation of block of web pages and removal of noisy blocks enable one to produce much better result for extracting useful information.

II. WEB CONTENT MINING

Web content mining identifies the useful information from the Web Contents/data/documents. However, such a data in its broader form has to be further narrowed down to useful information. In this section we begin with two main approaches of Web Content mining and define how it differs from Data Mining. The web content data consist of structured data such as data in the tables, unstructured data such as free texts, and semi-structured data such as HTML documents. The two main approaches in WCM are (1) Unstructured text mining approach and (2) Semi-Structured and Structured mining approach.

1. Unstructured Text Data Mining (Text Mining)

Web content data is much of unstructured text data. The research around applying data mining techniques to unstructured text is termed knowledge discovery in texts (KDT), or text data mining, or text mining. Hence one could consider text mining as an instance of Web content mining. To provide effectively exploitable results, preprocessing steps for any structured data is done by means of information extraction, text categorization, or applying NLP techniques.

2. Semi-Structured and Structured Data Mining

Structured data on the Web are often very important as they represent their host pages, due to this reason it is important and popular. Structured data is also easier to extract compared to unstructured texts. Semi-structured data is a point of convergence for the Web and database communities: the former deals with documents, the latter with data.

The initial step in data mining appears near is web content mining. It has been become complicated by the crawlers to find the well-organized data from deepness of web for extracting the required subject. Scanning the text also graphics so as to discover the applicable data is what all about content mining. To take out from the strength of web we need habitual tools for query able databases. The first step is structure mining that helps to cluster the web pages also more to take the scanning resulting into answers to the advanced level of application, therefore on the deep web it is promising to find out the results precedence to the search engines according to the highest application of keywords .It has also

in use the benefit of semi-structured character of web page, texted consist of the most functional prototype where text is followed by choice list with numeric values. Web content mining requires appearance of new and inventive applications and should also have its own individuality.

III. WEB CONTENT MINING TOOLS

A. Web Content Extractor

"Web Content Extractor" is software designed for web scraping, data mining, and data extraction. Web Content Extractor will permit users to mine the target data from a range of WebPages over the Internet.

Web Content Extractor can collect data from online stores, company directories, e-commerce web sites, economic web sites, shopping web sites, search engine outcomes, everything you can imagine that is going on the World Wide Web.

B. Web Info Extractor

"Web Information Extractor" is a very powerful tool used for web data mining and content mining, content investigation. It is able to extract structure or unstructured data from web page, alteration into local file or save to database, place to web server. No need to define difficult template rules, immediately browse to the web page you are interesting and hit it off what you wish for define the extraction job, and run it when you want, or allow it run automatically.

C. Web Text Extractor

"Web Text Extractor" is plan for extract text from web page and still control label in dialog simply. You can pull out and copy these texts with no select them. Mine text from web page, No text selection desired, Can mine unselect able text. Sort out transparent character as well as zero size character automatically. Mine text from still control, edit control and windows title. Handle extracted text for you.

D. Screen Scraper

Software that permits a PC to catch character-based data from a mainframe repeatedly presented in a green screen and it in an easier to recognize graphical user interface. Latest screen scrapers provide the information in HTML, thus it be able to access with a browser. Top producers include Mozart, lashpoint, Inc, and Intelligent Environments. An in-built recorder presents only click screen scraping.

Screen-scrapers allows mining the content from the web, like searching a database, SQL server or SQL database, which interfaces with the software, to achieve the content mining requirements. The programming languages like Java, .NET, PHP, Visual Basic and Active Server Pages (ASP) can also be used to access screen scraper.

E. Mozenda

Mozenda is software that permits commercial and nontechnical users to simply mine data across web pages. Mozenda now supports logins, paging throughout lists of

results, AJAX, frames, with other difficult web sites. Mined data can be accessed online, exported, as well as used throughout an API. Mozenda Data Extractor is an excellent tool that performs your scraper within the clouds. The circulated character of this web ripper works glowing for large amount scraping and listed and parallel web crop. Mozenda's service used for choosing items as well as appending harvest files fits well for grouping of data from various sources. The out coming export and distributing services (including email notifications) are wonderful characteristics of this screen scraper.

F. Automation Anywhere

Automation Anywhere is a web data extraction tool used for retrieving web data effortlessly, screen scrape from web pages or use it for web mining. The Intelligent automation software, used for automation and scheduling business process and IT tasks in easier way. Unique SMART automation Technology automates.

IV. DATA MINING TECHNIQUES USED IN WEB CONTENT MINING

Web content mining is the use of data mining techniques to automatically discover and extract information from web documents and services. Web content mining can therefore be effectively handled by various data mining techniques.

1. Association analysis: This method is used to find association/correlations of recurring pattern among item sets. This technique is used to discover association between groups of users with specific interests.
2. Clustering: This technique is used to group together items with similar characteristics. Clustering is used in creating user clusters based on navigational behavior which are derived from web logs. Clustering is also used to content by creating that are related in terms of their content. There are various clustering techniques as follows – Text based clustering, Partitioned clustering, hierarchical clustering, Graph based clustering, Neural Network based clustering.
3. Sequential Pattern Discovery: this technique combines the concept of association mining along with time sequences. This system uses server logs along with web content mining.
4. Segmentation : This technique is used to group users into various segments by using information from users profiles and past browsing/purchasing history.

CONCLUSION

There are many concepts regarding the World Wide Web. We tried to expose Web content mining, one of the categories of Web mining. The term Web Content mining refers to a technique that is use for extracting useful knowledge from web. We discuss different views to understand Web content mining, and given six tools that used in web content mining. We talk about the data mining techniques used in web content mining. Web content mining has been provided extremely functional in the business world. We observed that tools are user friendly and some of these tools seem to be applicable for data mining.

REFERENCES

- [1] Qingyu Zhang and Richard S. Segall, "Web Mining: A Survey of Current Research", Information Technology and Decision Making, **7(4)**, 683-720, 2008.
- [2] R Kosala, H Blockeel-ACM SIGKDD Explorations Newsletter, 2000.
- [3] Web Info Extractor Manual.
- [4] Ms Aparna Bulusy, "Web personalizing the web experience through data Mining", 2nd international conference on Emerging Trends in computer science, communication and information technology, Feb -2015
- [5] Tripurari Pujan pratap singh and K.K. pandey, "HIT: Web content mining tool" International Journal of Electronics Communication and Computer Engineering Volume 3, Issue 6, ISSN (Online): 2249-071X, ISSN (Print): 2278-4209
- [6] C. Lakshmi Devasena Et AL. International Journal on Computer Science and Engineering IJCSE), ISSN : 0975-3397, 3 Mar 2011, Vol. 3 No. 1155-1167
- [5] www.screen-scraper.com
- [6] Web Content Extractor help.
- [7] Automation Anywhere 5.5 help
- [8] V. Bharanipriya and V. Kamakshi Prasad, "Web content mining tools: A comparative study," International Journal of Information Technology and Knowledge Management, vol. 4, no. 1, pp. 211-215, Jan 2011.