

# Extract Structured Data from Heterogeneous Web Pages

M.D. Nirmal<sup>1</sup>, Shital B. Jadhav<sup>2</sup>, Nilam V. Dhumal<sup>3</sup>, Rutuja V. Kapadnis<sup>4</sup>, Priyanka H. Thakare<sup>5</sup>

Computer Engineering, Pravara Rural Engineering College, Loni, Maharashtra, INDIA<sup>12345</sup>

**Abstract:** Information Extractor is a powerful tool for web data mining and data crawling. Data from web pages. Reform into local file or save to database, post to web server. No need to the web page you are interesting and click what you want to define the extraction task, and run it as you want, or let it run automatically. Data Extraction is act or process of retrieving data out of data sources for further data processing or storage. The import into the intermediate extracting system is thus usually followed by data transformation and possibly the addition of metadata prior to export to another stage in the data workflow. We formally define structured data, the kind of data that we are hoping to extract from the web pages Structured Data is any set of data values conforming to a common type. The *Basic Type*, denoted by, represents a string of *tokens*. A token is some basic unit of text.

**Keyword:** Main Content Extraction (MCE), HTML&XML, Data Mining, data crawling, Data scrapping

## I. INTRODUCTION

World Wide Web is now a medium by which people all around the world can use this and gather the information of all Types of data. The various sites web pages are generated dynamically contain semi structured information also. This information is mainly, advertisements, copyright statements, privacy statements, logos, table of contents, navigational panel, footers and headers come under noisy content. The irrelevant data will be present in the web page information. User want to need the important or usable data into that pages. At that time our project is use to extract those data that user want from web page.

The large amount of information on web is stored in backend databases which are not indexed by traditional search engines. Such databases are referred to as Hidden web databases and extraction of this hidden web content is a potential research area as the pages are dynamically created through search query engine. However, direct query through this search engine is laborious way to search. Hence, there has been increased interest in retrieval and integration of hidden web data with a view to give high quality information to the web.

One of the important task in information retrieval science is to extract useful content from web pages. The problem is web pages (usually HTML) content some information (text) with tons of information additional texts-navigation structures, advertising etc. Form point of view of information retrieval tool each html page has main (useful) content and helpful information that is good when viewing web, but not when extracting the data. The task is to filter useful content from HTML without knowing of structure of the page. It is not too difficult to get useful content when structure of HTML page known. For example, with using web scrapper. But it is not possible to create templet or rage expressions for

every page in the web. There is simple solution for this task. As it is known, HTML page can be presents as tree view.

The rapid development of the internet and web publishing techniques create numerous information sources published as HTML pages on World Wide Web. However, there is lot of redundant and irrelevant information also on web pages. Navigation panels, Table of content (TOC), advertisements, etc. Generally user want to interested only main contents of the web so we remove unwanted information as per requirement. Then for this purpose various technics are used like data mining, data crawling, topic distillation etc. This paper discusses various approaches for extracting informative content from web pages and a new approach for content extraction from web pages using word to leaf ratio and density of link.

Web Miner, Commercial software for extraction specific information, images and files from websites. E-commerce product and pricing database that obtains its data through information extraction from thousands of online retailers. There are number of web data Scrapping Tools Available on the internet. Most of the tool are offered by data mining companies. For the last decade the internet has revolutionized the way project handle information. Project have taken the advantage of the huge amount of information find on internet. That will need for a project to search numerized and important web site for content and information that is relevant to the task that is needed to be performed. Data scrapping service to research companies reduces greatly the cost in incurred in business. Data mining company's employer expert in data mining. The experts are extract data from different websites and also compare data extracted from another web sites. You can get a lots of information from many web site in the shortest time possible.

Even you are using many application, web data scrapping, data mining tool offered will be useful in the getting

right information at the right time. Web Scrapping is the complex process. This explain why it is perform by experts as specialized company. The data which is extracted from internet has wide application in various industries and if the data is proceed it can be use be greatly. Most data mining specialized in web destination, data mining and other data management services. Web data extraction this involve the process of extracting web pages of text-based mark-up language such as HTML. Which surely contain the welch of useful information.

**The user needs to follow the following steps.**

1. He/ She has to discover the URLs of the hidden Websites.
2. Visits homepages of these web sites.
3. Send queries through HTML forms.
4. Extract the relevant information from the result web pages.
5. Compare or integrate the results from multiple sources.

**II. LITERATURE SURVEY**

Information extraction dates back to the late 1970s in the early days of NLP. Data scrapping service to research companies reduces greatly the cost in incurred in business. Data mining company’s employer expert in data mining and web scrapping. The experts are able to extract data from different websites and also compare data extracted from different web sites. You can get a lots of information from many web site in the shortest time possible.

Janet White try Semantria’s API. 1st tried it using their free trial. He was able to pull names (they call them entities) and themes that reoccurred in my content. By creating queries you can also tag the common topics present in your group of texts. Then Antonio Matarranz tell Textalytics (/Meaning as a Service) is a cloud-based semantic API that offers a topic Extraction service (entities, concepts).

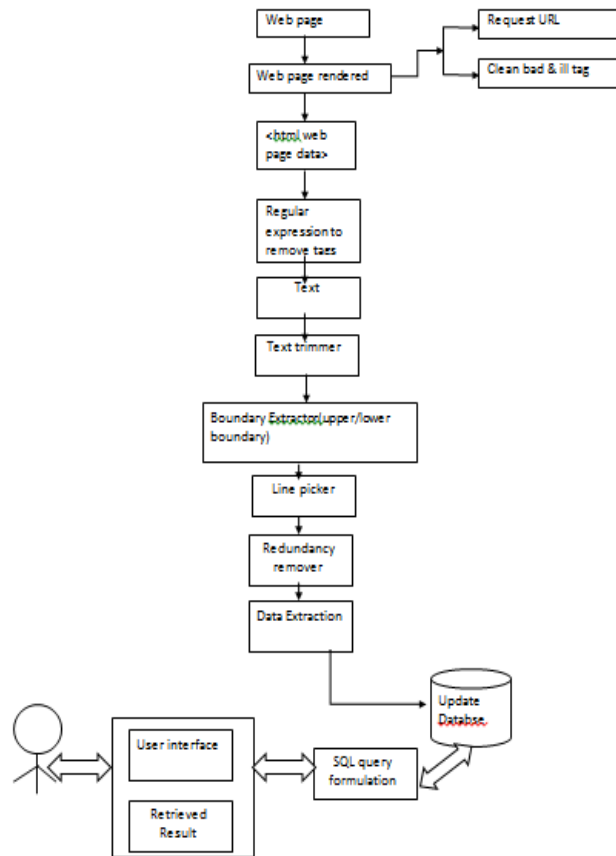
In addition, you can tag your text with theme categories, feature-level sentiment, etc.

Text Trimmer1 takes plaintext copy of the extracted data and does the same thing that Comma Remover does but alongside it also removes white spaces. It first replaces all whitespaces with commas and then removes all commas generated and text displayed on web page). The approach used in this paper for this module is static that is user has to select a relevant data with the help of a mouse.

An easy way to comply with the conference paper formatting requirements is to use this document as a template and simply type your text into it.

**III. PPROJECT WORK**

HTML document is main Content that we want to prefer or extracting data with the help of that Pages. Following Architecture Show those things that easily understandable for this project.



**Web Page:** A web page is a document that is suitable for the World Wide Web and the web browser. Web browser display a web pages on a monitor device. The web page is what displays, but the term also refers to a computer file, usually written I html or comparable mark-up language. First enter into web page then go into web page render.

**Web Page Render:** You type an URL into address bar in your preferred browser. The browser parse the URL to find the protocol, host, port, and path. The browser receives the response and parse the HTML (which with 95% probability is broken) in the response. In that clean the ill tag which was created at the time of extract data. Then after that html web page data come. That contain tots of tag which present on HTML page of given site.

**Regular Expression to Remove Task:**

We can use regular expression for regulate HTML to parse the HTML code. HTML is not a regular language and hence can’t be 100% correctly parsed with a regex. This is just one of many problems you will run into. The best approach is use an HTML parser to do this for you.

**TEXT:** HTML also defines special elements, for defining text with a special meaning. HTML uses elements like <b> and <i> for formatting output like bold or italic text. Formatting element were designed to display special types of text like Important text, Emphasized text, Marked text, Small text, Deleted text, Inserted text, etc.

**TEXT TRIMMER:** When we need to rearrange a text list you might be faced with borrowing operations. The Text Trimmer program may help you addressing two cases: Block operations and Mass operations. For important the text trimmer is use to remove white spaces (blank space) into the HTML document.

**BOUNDARY EXTRACTOR:** Boundary Extractor discover of information from unstructured or semi structured. We capture the structure of a document as a tree of nested HTML tags. Header and Footer removed with the help of Boundary extractor.

**LINE PICKER:** There are multiple tags Present before the boundary extractor is work after line picker remove those unwanted tag & fed them as Less Copy of tag is call Line Picker.

**REDUNDANCY REMOVER:** This algorithm checks multiple occurrences of commas with a single comma and marks it as a field separator. This process also takes care of the fact that there should not be any leading and trailing commas so it detects such commas and removes them from the file

**DATA EXTRACTOR:** Web data Extractor it uses Regular Expression to find Extract and scrape internet data quickly and easily. It is very Flexible allowing you to extract both data and complex data structures like HTML tables. The data Will Extracted over here after that redundancy remove is completed.

**UPDATE DATABASE:** Before Creating Database First Blank Database Will Created Then after extracting the sorted data (pured data) it will create title for the content of data. The separated In the Database by filling the information purpose. Finally Database Will be created that user need.

**SQL QUERY FORMULATION :** We want to perform various operation in our New Database then the sql query Evaluation is used to perform some insert, delete, update, alter, drop, etc. operation.

#### IV. ALGORITHM

1. Construct DOM (Document Object Model) tree of web page.

2. Remove unwanted nodes which contain title, script, and head, Meta, style, no script, link, select, #comment and the nodes which have no words and are not visible.

3. Compute word to leaf ratio (WL) as in Eq. 1.

$WL(no) = w(no) / l(no)$  (1) Where  $w(n)$  = number of words in the node no

$l(no)$  = number of leaves in the subtree of node no

4. Obtain starting node set which contains higher density of text. Initial node set IN is defined as those nodes which satisfy the condition in Eq. 2.

$WL(no) \geq \sqrt{\max WL \times WL(\text{root})}$  (2)

5. Compute text link ratio (TL) i.e. the ratio of the length of the text and the number of links of node n in IN.

6. Compute link text ratio (TLT) i.e. the ratio of the length of the text and the length of the link text of node no in IN.

7. Calculate weight of each node  $W(no)$  as in Eq. 3

$W(no) = (TL(no) + TLT(no)) / b$  (3)

Where b is normalizing factor.

8. Relative position of node no is defined in Eq. 4, Eq. 5, Eq. 6 and Eq.7.

$R(no) = WL(no) \times \max w(no), \Sigma R(no1)$

$No1 \in \text{Children}(no)$  (4)

Where  $w(no) = rpos(no) \times rWL(n)$  if  $n \in IN$  (5)

0 if  $no \in IN$

Where  $rpos(no) = 1 - (id(no) - mnid) / mxid - mnid$  (6)

$rWL(no) = (WL(no) - mnWL) / (mxWL - mnWL)$  (7)

Where mnid, mxid are the minimum and maximum values of identifiers in IN, mnWL, mxWL are the minimum and maximum WLR values from step 3.

Relevance includes those nodes which have higher density of text. If two nodes have same  $R(no)$  then the node with lower identifier is selected.

9. Select the node which have highest values of  $W(no)$

And  $R(no)$  in proportion as,  $0.7 \times R(no) + 0.3 \times W(no)$ .

10. Selected node has the required textual information.

#### V. CONCLUSION

Now a days web pages are unstructured and its number is growing at a very fast rate for these web pages informative content extraction is very important. Unstructured web pages contains repetitive and irrelevant information so it reduce the performance and consumes lot of time. Therefore to avoid those issues, an approach for extract structured data from web pages has been implemented in this proposed work. It is useful for human users because it will give required information in time efficient manner. It also used as pre-processing stage for applications that need to extract only the main content of web pages (for example robots, indexers, crawlers etc.) to prevent the processing of useless, irrelevant and noisy information.

#### VI. REFERENCES

- [1] [1] Bing Liu, Robert Grossman, and Yanhong Zhai. Mining data records in web pages. In KDD '03: Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining, pages 601–606, New York, NY, USA, 2003. ACM Press.
- [2] [2] Nodules, A., Zeros, P., Cho, J. Downloading Textual Hidden Web Content Through Keyword Queries. In Proceedings of the 5th ACM/IEEE Joint Conference on Digital Libraries.
- [3] [3] Ji Ma; DerongShen; TieZhengNie DESP: An Automatic Data Extractor on Deep Web Pages Web Information Systems and Applications Conference (WISA), 2010 7th Publication Year: 2010, Page(s): 132 - 136
- [4] [4]Anuradha, A.K.Sharma, "A Novel Approach for Automatic Detection and Unification of Web Search Query Interfaces using Domain Ontology" selected in International Journal of Information Technology and knowledge management(IJITKM), August 2009.
- [5] [5] JianQiu, Feng Shao, MishaZatsman, JayavelShanmugasundaram, Index Structures for Querying the Deep Web, Workshop on the Web and Databases (WebDB), 2003, 79-86
- [6] [6] S. Raghavan and H. Garcia-Molina. Crawling the Hidden Web. In Proceedings of VLDB, pages 129–138, 2001.
- [7] [7] Nodules, A., Zeros, P., Cho, J. Downloading Textual Hidden Web Content Through Keyword Queries. In Proceedings of the 5th ACM/IEEE Joint Conference on Digital Libraries (JCDL05). 2005.