

# Extraction of Information from Web Page Using Content Mining Approach

Yogaraj Manickam<sup>1</sup>, Rajalakshmi<sup>2</sup>, Sambath Uma<sup>3</sup>

<sup>1,2,3</sup> Dept. of CSE, Park College of Engineering And Technology, Tamil nadu, India.

## To Cite this Article

YOGARAJ MANICKAM<sup>1</sup>, RAJALAKSHMI<sup>2</sup>, SAMBATH UMA<sup>3</sup>, "Extraction of Information from Web Page Using Content Mining Approach", International Journal of Scientific Research in Engineering & Technology, Volume 02, Issue 02, March-April 2022 PP: 13-15.

**Abstract:** Today internet has made the life of human dependent on it. Almost everything and anything can be searched on net. The quick improvement of World Wide Web has been gigantic of late. With the large amount of information on the Internet, website pages have been the normal wellspring of information recuperation and data mining technology, for instance, business web crawlers, web mining applications. In any case, the site page as the chief source of data contains many parts which are not likewise critical. Other than the major things, a site page in like manner incorporates of noisy parts that can spoil the display of information recuperation applications. In this manner cleaning the site pages before mining becomes essential for additional fostering the mining results. In our work, we bases on recognizing and wiping out local noises in site pages to deal with the show of mining. The information contained in these non-content blocks can distract the client and moreover hurt web mining. So, it is basic to confine the instructive fundamental substance blocks from non-valuable blocks. Along these lines, we propose a system that kill different upheaval plans from any page. There are two steps, Web Page Segmentation and Informative Content Extraction, are expected to have been finished for Web Informative Content Extraction. We will inspect the site page and by using methods and estimation we eliminate topic information requested by user.

**Watchwords** — Web Mining, Web Content Extraction, DOM Tree, Information retrieval, HTML Parser

## I. INTRODUCTION

With This loud information makes extraction of Web content bleak. Various strategies are accessible for web content extraction. The utilization of data mining techniques to thus find and to remove data from Web data, including Web documents, hyperlinks between documents, usage logs of Websites, etc, is called Web mining. Some of the data mining procedures applied in Web mining are alliance rule mining, batching, portrayal, ordinary thing set. Some of the sub tasks of Web mining are finding of relevant resource, assurance of information and pre-taking care of, theory and analysis. Web content digging is used for isolating significant information from Web pages. Site page content can be structured, unstructured and semi-coordinated. Coordinated Web page data are easy to eliminate when differentiated and unstructured and semi-coordinated data. Web Content Extractor consistently eliminates a whole Web page including joins, header, footer, essential substance and advertisement. During the extraction unwanted data like associations, header, footer and promotion are treated as loud information. To discard the riotous information and concentrate the significant information is a troublesome issue. Various methods were proposed for taking out rowdy information. Exactly when a client question the web using the web search instrument like Google, Yahoo, Alta Vista etc, and the search engine returns thousands of links related to the keyword searched. Now if the first link given by the user has only two lines related to the user query & rest all is uncluttered material then one needsto extract only those two lines and not rest of the things. The current study focuses only on the core content of the web page i.e. the content related to query asked by the user. The title of the webpage page, Pop up advancements, Flashy promotions, menus, silly pictures and associations are not relevant for a user querying the system for educational purposes.

## II. RELATED WORK

This study is proposed to deal with the problem of intra-page redundancy that causes search engines to index redundant contents and recuperate non-critical results. The issue moreover impacts Web diggers since they eliminate plans from the whole document rather than the illuminating substance. Along these lines, we outline examinations of the two fields. In the rest of the paper, for better understanding, we use information recuperation (IR) systems to imply web crawlers and information extraction (IE) structures to mean Web or text diggers. Various IR structures have been done to normally collect, cycle, record, and separate the Web documents for serving clients information needs. It also parses things in the page considering HTML or other increment language like XML. The former called text mining. jsoup is designed to deal with all varieties of HTML found in the wild; from pristine and validating, to invalid tag-soup; jsoup will create a sensible parse tree.

### III. EXAMINATION

Loud substance makes the issue of information gathering from site pages significantly more earnestly. Website pages consistently contain non-edifying substance, disturbances that could antagonistically impact the introduction of Web Mining. Exactly when a client request the web using the web crawler like Google, Yahoo, AltaVista, etc, and the web search instrument returns enormous number of associations associated with the keyword searched. By and by accepting the primary association given by the client has basically two lines associated with the client question and rest everything is tidied up material then one necessities to remove simply those two lines and not rest of the things. Considering that a tremendous proportion of world's information re-sides in web pages, it is becoming increasingly important to analyze and mine information from web pages.

#### A. Identifying Articles

The first step, determining whether a page contains an article, is a document classification problem. Our evaluation implicitly expects that such a classifier is given, since all our testing models contain articles. No such assumption that is made in training, regardless, and the semi-normally delivered getting ready data may erroneously contain non-articles.

#### B. Cleaning Extracted Blocks

Most outrageous delayed consequence division keeps an eye on the second supporter of the issue, recognizing the block of HTML containing the article message, yet further cleaning may be supposed to take out unessential words from embedded advancements, game plans of associations with other stories, images with captions, eye-getting accentuations of proclamations appearing in the article, etc. That is the thing a key insight is, the point at which the article's HTML block has been recognized, dispensing with whatever "trash" remains ends up being significantly less troublesome. Undoubtedly, after inspecting our evaluation sets we found this ought to be conceivable using several guidelines with little botch. Starting with an isolated block of HTML that starts at the first word and ends at the last word of the main text of the article:

This heuristic capabilities commendably considering the way that the text of most reports only from time to time contains interfaces, and inserts are wanted to either show you something (for instance a picture) or motivate you to make a beeline for some spot (for instance an association). The several slips up we saw were mostly trivial (like leaving in "Business"), but one article would lose most of its text since it contained links and was inside a <DIV>. Outside the news space we will not really be so fortunate: reference book articles, for example, tend to contain numerous associations, and would require a more complicated technique. All things being equal, in any case, cleaning the removed block is far less critical than separating it precisely, since site experts consistently place the greatest proportion of trash text, such as navigation bars or client comments, around the article, not inside it, paying little brain to space. This is displayed by our very strong cross-region execution cleaning general articles from the Clean Eval task paying little heed to advancing no endeavor to clean the blocks selected by MSS.

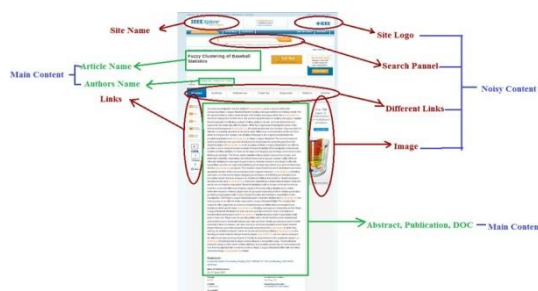


Fig.1 webpage of IEEE Explore article

### IV. PROPOSED WORK

Proposed approach centers around site pages where the secret information is unstructured text. The technique used for information extraction is applied on entire web pages, whereas they actually seek information only from primary content blocks of the web pages. The user specifies his required information to the system.

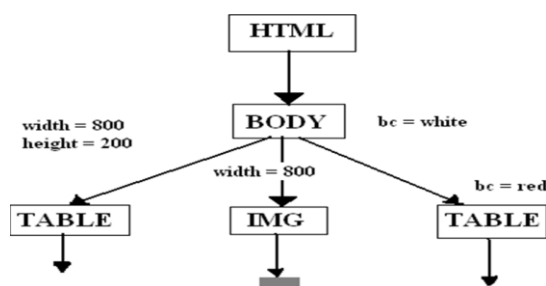


Fig2: DOM Tree

```
<HTML>
  <BODY bgcolor=WHITE>
    <TABLE width=800 height=200>
      ...
    </TABLE>
    <IMG src="image.gif" width=800>
    <TABLE bgcolor=RED>
      ...
    </TABLE>
  </BODY>
</HTML>
```

Fig3:HTMLcode

#### A. Jsoup

jsoup is a Java library for working with certified HTML. It gives an incredibly supportive API to isolating and manipulating data, using the best of DOM, CSS, and jquery-like methods. jsoup implements the WHATWG HTML5 specification, and parses HTML to the same DOM as modern browsers do.

- scrape and parse HTML from a URL, file, or string
- find and extract data, using DOM traversal or CSS selectors
- manipulate the HTML elements, attributes, and text
- clean user-submitted content against a safe white-list, to prevent XSS attacks
- output tidy HTML

#### V. CLOSES

This paper proposed an novel task for finding local noise in web pages. Using DOM tree approach contents of the web pages are extracted by filtering through non enlightening substance. With the Document Object Model, engineers can manufacture documents, navigate their plan, and add, change, or eradicate parts and content. With this features it turns out to be more direct to isolate the useful content from a large number of web pages. In future this approach will be used in information retrieval, automatic text categorization, topic tracking, machine translation, abstract summary. It can provide conceptual view of document collections and has important applications in the real world.

#### REFERENCES

1. D. Cai, S. Yu, J.-R. Wen, and W.-Y. Ma. *Vips: vision-based page segmentation algorithm*. Technical report, Microsoft Research, 2003.
2. A. H. F. Laender, B. A. Ribeiro-Neto, A. S. da Silva, and J. S. Teixeira. *A brief survey of web data extraction tools*. *SIGMOD Rec.*, 31(2):84-93, 2002.
3. Y. Yesilada, — *Web Page Segmentation: A Review*, leMINETechnicalReportDeliverable0(D0), 2011.
4. [6]. Y. Yesilada, — *Heuristics for Visual Elements of Web Pages*, leMINETechnicalReportDeliverable1(D1), 2011.
5. Zhao Xin-xin, Suo Hong-guang, Liu Yu-shu. *Web Content Information Extraction Method Based on Tag Window*. *Application Research of Computers*. 2007, 24(3). -144-145, 180.
6. Pan Donghua, Qiu Shaogang. *Web Page Content Extraction Method Based on Link Density and Statistic*. *The 4Th International Conference*.
7. A. F. R. Rahman, H. Alam and R. Hartono "Content Extraction from HTML Documents"
8. Wolfgang Reichl, Bob Carpenter, Jennifer Chu-Carroll, Wu Chou "Language Modeling for Content Extraction in Human-Computer Dialogues".