

Extraction of Textual Contents From News Web Pages Using Pattern Matching

Pramoh Jain¹, AryanJain², Avnish Jain³, Bhagyesh K Shah⁴, Yashpal Gupta S⁵

^{1,2,3,4}Students, Computer Science and Engineering Department

⁵Guide, Assistant Professor, Information Science and Engineering Department

Bahubali College of Engineering, Shravanabelagola

ABSTRACT: This Research work deals with extraction of textual content from web pages and also deals with ranking of news channel based on sentiment analyzer. Web pages besides textual content consist of other elements, such as banners, navigational elements, copyright information, external links, these are called noisy content. Most of the information available on web pages is either represented in XML, or HTML, or XHTML format that mostly contains semi-structured text documents. This document does not discriminate between the text and the schema, and the amount of structure used to represent the text depends on the purpose. No semantic is applied to semi-structured documents. This requires extracting core contents of text document to analyze words or sentences for retrieving relevant information. Although there are many existing methods that formulate the actual content identification problem as a DOM tree node selection problem, each one has some sort of lacunae. Here we proposed an approach based on pattern matching technique. This technique uses simple heuristic for extraction of core contents from web pages which are mostly semi-structured in nature. It requires visiting the appropriate news web site using their URL, accessing the links related to each news page of specified category, extracting the data including metadata from each of these news web pages. The approach uses devised algorithm that applies regular expressions (regexes) to identify the correct pattern for extracting the actual text contents from these news documents. Proposed approach deals with news web pages of any size and extracts core contents with efficiency and high accuracy.

Keywords: Pattern matching, Information extraction, Document Object Module, Sentiment Analyzer tags.

I. INTRODUCTION

Rapid growth of information is taking place on the web and huge amount of information is available in the form of online repositories, such as web pages, web archives, digital libraries, etc. Heterogeneous sources of online repositories mostly contain textual data available either in a semi-structured or an unstructured form. Extracting relevant information from these heterogeneous sources of semi-structured or unstructured form is useful for many applications, such as document identification, text categorization, summarization, clustering, topic tracking, etc. Several techniques have been evolved in view to provide an efficient access to relevant information from these online repositories. Extracting relevant information from unstructured text data results in development of breakthrough techniques based on Text mining.

Many techniques adapted the method of extracting core contents from the web pages pre-process the data, and then applying the process of Information Retrieval (IR) or Information Extraction (IE) for retrieving the relevant information. Most of these techniques are not widely applicable as IR based ones, because they either need human intervention, and/or the quality of extraction is very low. Thus there is a scope to develop an efficient and effective method for extraction of high quality contents from semi-structured data available on the web.

In this Research work, we present a template independent approach based on pattern matching for extraction of core contents from news web pages. This approach extracts the data from each news web pages that are organized as well formatted XHTML documents. The devised algorithm is used that applies the process of pattern matching for identification of title element and the body of the news article, extract the core contents by filtering out the noise, and stores these contents into plain text form

II. RELATED WORK

Web page extraction techniques/approaches are mostly based on either generating wrappers [1], or tag based or tree based approaches [5][6][7][8], or using the techniques of machine learning[2][9], or natural language processing[1][3]. These approaches are either domain-dependent, and/or need human intervention, or require large amount of dataset for training purpose [10]. Sandeep Sirsat [10] has discussed several techniques and approaches useful for extracting core content from semi-structured text documents, and explores merits and demerits of each of these techniques. Many techniques use Tag-based approach that uses HTML parser to convert each web page into DOM tree. It then applies traditional pattern reduction techniques in order to find a template. Most of these techniques are template dependent and include substantial features of the DOM tree.

Many researchers utilizes the fact that web pages are naturally semi-structured in nature and utilizes the feature of DOM tree by representing it as a labeled ordered routed tree. These techniques mostly rely on analysing the structure of the target pages [13][4]. Some techniques uses the concept of tree did distance to evaluate the structural similarity between pages [5][6][7][8]. Mohsen Asfia, Mir Mohsen Pedram, Amir Masoud Rahmani [13] introduced an algorithm to simulate a web page and to find the main content block position in the page. The proposed method is tag independent and has two phases to accomplish the extraction of main content. The method could not have any learning phases and could find informative content on any random input detailed web page. Furthermore it could detect and eliminate comments from the extracted content. Y. Lou, Y. Z. Yuan, Zhang [4] proposed a method for extraction of website information based on DOM to improve the searching efficiency. This method preserves the theme information and to filter out the noisy content that the users are not interested in. It presumes that subject content extraction of pages has been essential for Web information pre-treatment link. It can reduce the browsing time, promote the speed of user accessing to information, improve efficiency and enhance the usability of the web, with extracting the topic information by DOM model.

III. Identification & Extraction of Title elements:

The title of the news document is represented by the title element, which is mostly enclosed in title tag. The

module locates the start and end of title tag and then extracts the text content in between.

```
<!DOCTYPE html PUBLIC "-//W3C//DTD
XHTML 1.0 Strict//EN"
"http://www.w3.org/TR/xhtml1/DTD/xhtml1-
strict.dtd">
<html xmlns="http://www.w3.org/1999/xhtml"
xml:lang="en" lang="
xmlns:fb="http://ogp.me/ns/fb#">
<head>
<meta http-equiv="Content-Type"
content="text/html; charset=UTF-8
<title>
Gambhir- man on a mission - The Hindu
</title>
<link rel="shortcut icon" type="image/x-icon"
href="http://www.thehindu.com/favicon.ico" />
<link rel="icon" type="image/ico"
href="http://www.thehindu.com/favicon.ico" />
```

The full text story which represents the actual contents of the news web page is placed either in <div> or <p> tags followed by class attribute value which may differ from different e-newsResearch work to different e-newsResearch work.

```
<div class="article-body">
<p class="body">GautamGambhir appeared a man on a mission
after his omission from the Indian Test squad.
</p><pclass="body">He was light on his feet and weighty with
his
```

Most e-newsResearch works contain noisy contents such as images, advertisements, references to named entities, etc. present in the body of the full text story it may also include noisy elements such as
 tags that represent line break, <p> tags, Anchor tags, apostrophes, etc. The algorithm filters out the noisy contents and extracts high quality contents with efficiency and accuracy. For this, substitute regexes are defined for filtration of such noise from the textual content of the story. This transforms the content of news web pages automatically into the plain text form, herein after, called as text document

WEB CRAWLER

First enter the website information and it will search the information using Google URL. The Google result initially it will validate the Google URL if it is correct then it will extract the data from the correct Google URL after that it will store the data into the repository, only after validation of the Google URL the data can be stored. The Google result process is completed.

DATA CLEANING

The Data Cleaning algorithm is responsible for removal of stop words. These are the set of words which do not have any specific meaning. Data Cleaning is used for removing the stop words from each of the news and clean them. After the data cleaning process is completed the clean data can be represented as a set (CleanId, Clean Data, Description). Clean Id is the unique Id associated with the news, Clean Data is the clean data after removal of clean data.

TOKENIZATION

Tokenization is a process of converting the clean data into a set of words known as tokens. Each of the token can be represented as Token ID, Token Name and Application ID.

Filtration of Noisy Content:

```

Read data from stored news document;

Split the data on whitespaces & store them in an array
    @words[];

Identify the pattern containing name of newsResearch
work and replace it
    by blank.

Foreach $word in an array @words[] loop
    Part 1: Extraction of TITLE element.
    If $word includes pattern "<title>" then
        Set $i = 1;
        Replace $word by blank;
    Elseif $word contains pattern "</title>" then
        Set $i = 0;
        Replace $word by blank;
    endif
    If ($i) then
        Extract ($word);
    Check $word to filter noise and store it;
    endif
    Part 2: Extraction of Core contents.
    If $word includes pattern "<div class = $regex>" Or
        "<p class
        = $regex>" then
        Set $j = 1;
        Replace $word by blank;
    Elseif $word contains pattern "</div>" then
        Set $j = 0;

```

F POLARITY N

The clean data is examined carefully for positive, negative and neutral keywords. Separately different-different words are categorized as positive keyword, negative keyword and neutral keywords. Then after retrieving different category of keywords positive polarity for positive keywords, negative polarity for negative keywords and neutral polarity for neutral keywords. Finally a polarity matrix is formed and based on feature vector the ranking is done.

RECOMMENDATIONS BASED ON SENTIMENT ANALYSIS USING IDFT COMPUTATION

The Sentiment Analyzer Agent is used to rate the News into sentiments like positive, negative or neutral sentiments by matching a set of statements by using dynamic statement feeder and sentiment type. The sentiment analyzer agent which consist of Sentiment Repository which consist of sentiments i.e statements of sentiments and its type. The sentiment rater is used to extract the sentiment from News and rate it as positive, negative and neutral.

V. CONCLUSION

In this Research work, we proposed an algorithm for extracting the core contents of news web pages using pattern matching approach that transforms the contents of news web pages automatically in to plain text form. This approach does not exploit the features of DOM structure. It is Template independent and could not dependent on any tag type. It has been observed that the accuracy of extraction of core contents after transforming the documents available in XHTML format to plain text form is efficiently high. This approach deals with news web pages of any size and extracts core contents with efficiency and high accuracy. The algorithm extracts high quality contents with efficiency and accuracy even from short news web pages.

REFERENCES:

- [1] Chia-Hui Chang, Mohammed Kayed, MohebRamzyGirgis, Khaled Shaalan - A Survey of Web Information Extraction Systems , IEEE TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING, T Vol.18, no.10,pp. 1411-1428, Oct. 2006.
- [2] J. Prasad and A. Paepcke, "CoreEx: content extraction from online news articles," in CIKM '08: Proceeding of the 17th ACM conference on Information and knowledge management. New York, NY, USA: ACM,2008, pp. 1391-1392.
- [3] Tanveer Siddiqui, U. S. Tiwari, Natural Language Processing and Information Retrieval, Oxford University Press Pages 67-68.
- [4] Y. Lou, Y. Z. Yuan, Zhang, Website information extraction based on DOM-model, Proceedings of the 2nd International

Symposium on Computer, Communication, Control and Automation (ISCCCA-13), Atlantis Press, Paris, France, 2013, pp 792-795.18

[5] QiuJun LAN- Extraction of News Content for TextMining Based on Edit Distance Journal of Computational Information Systems 6:11 (2010) 3761-3777, November, 2010.

[6] Davi de Castro Reis, Paulo B. Golgher, Alberto H. F.Laender, Altigran S. da Silva, Automatic Web News Extraction Using Tree Edit Distance, ACM WWW2004, NewYork, USA, May17-22,2004.

[7] W. Chen. New algorithm for ordered tree-to-tree correction problem. Journal of Algorithms, 40:135- 158, 2001.

[8] K. Zhang, R. Statman, and D. Shasha. - On the editing distance between unordered labeled trees. Information Processing Letters, 42(3):133-139, 1992.

[9] S. Wu, J. Liu, and J. Fan, Automatic Web Content Extraction by Combination of Learning and Grouping, ACM WWW (IW3C2) May 18-22, 2015, Florence, Italy.

[10] Sandeep Sirsat, "Extracting Core Contents from WebPages", International Journal of Engineering Trends and Technology (IJETT) ISSN: 2231-5381 - Volume8 Number 9, pp 484- 489, February 2014,