



Functional Architecture of integrated framework for Facet-based Data Collection and Analysis

Dr Kanaka Durga Returi¹, Dr. Vaka Murali Mohan²
Professor^{1,2},

Department of CSE Engineering,
Malla Reddy College of Engineering for Woman,
Maisammaguda(V), Gundlapochampally, Medchal (M&Dist),
Hyderabad.-500100, Telangana, India.

Abstract — We present in this paper an integrated framework for collection and analysis of Facet-based text data. The integrated framework consists of four components: (1) user interface, (2) web crawler, (3) data analyzer, and (4) database (DB). User interface is used to set input Facet and option values for web crawling and text data analysis using a graphical user interface (GUI). In fact, it offers outcomes of research by data visualization. The web crawler collects text data from articles posted on the web based on input Facets. The data analyzer classifies papers in "relevant articles" (i.e., word sets to be included on such posts) and "nonrelevant articles" with predefined information. It then analyzes the text data of the relevant articles and visualizes the results of the data analysis. Ultimately, the DB holds the generated text information, the predefined user-defined expertise and the outcomes of data analysis and data visualization. We verify the feasibility of an integrated framework by means of proof of concept (PoC) prototyping. The experimental results show that the implemented prototype reliably collects and analyzes the text data of the articles.

Keywords— Data Analysis, Integrated Framework, Intelligent Service, Text Data Collection, Web Crawling.

I. INTRODUCTION

Intelligent systems have recently received significant interest from both academia and industry, for instance media remedy and choice research and recommendation.(1–3) Such resources often use text-data from papers posted on the web to gather information that people require. In general, data collection and analysis are the most important features of the Web system. The data can be captured, stored and processed using a "internet sensor" which is a special type of network-centered infrastructure. Therefore, web crawling to collect text data and Data Analysis to analyze collected text data are widely considered as key enablers of sensor web for such intelligent services.

To date, some current research projects have attempted to use open source programming languages such as R, Python, and Scala to incorporate such functionality(5–7) Though, most have never used an automated web crawling or big data analytic framework. This is, in most of the existing studies, to separate the functionalities of web crawling and big data analytics. In order to enable smart services to be seamless, an integrated architecture of different functionality (e.g. web crawling, data analysis and user application) must be designed (10) and is therefore subject to unpredictable

delays in a smart targeted service, since its feasibility highly depends on the developer's expertise.

In this article, they suggest an automated web-crawling platform and data analysis to allow smooth, intelligent services to collect and interpret Facet-based text content. The framework proposed consists of the following four components: (1) user interface, (2) web crawler, (3). These components interact to exchange data. The user interface component helps users to set a graphical user interface (GUI) for input Facet and detailed option values for web crawling and text data analysis. In comparison, different results for visualization of content, such as word clouds or word intensity charts, were generated based on results from the study of text information. The component of the Web crawler collects text data from web-based articles and provides data for the analysis through the storage of data from the collected text on the DB. The component of the data analyzer performs data pre-processing and analysis using data sets from the component of the Web crawler. The identification of objects is carried out for information preprocessing.

Throughout general, the papers will be categorized under predefined information (ie a set of words to be included) as "relevant articles" or "nonrelevant objects." Relevant articles are collected articles closely related to the subject the user is looking for, but non-relevant articles are not very closely related to the theme. Two steps are followed by data analysis. The first step is to extract words from the texts of corresponding articles, which consist of three or more characters. The latter is a process to filter words that eliminates unnecessary words. The data analysis results are displayed in word clouds and word frequency charts. In the last study, DB consists of three DBs: DB, DB and DB. The DB portion consists of three DBs. Every DB stores the text data collected, predefined knowledge and data analysis and visualization results. Through proof of concept (PoC) prototyping we check that the unified system is feasible. The web crawler and data analysis is performed by means of open-source R packages and the interface is implemented using Java Swing frames (11,12). The results from the

experiments show that the integrated frame offers the functionality of web crawling and text data analysis reliably.

II. FUNCTIONAL ARCHITECTURE OF INTEGRATED FRAMEWORK

Figure 1 displays a user interface, internet crawler, information analyser and DB element operational structure in the suggested unified system. There are three logical elements in the user's interface component: input panel, display panel and a predefined information panel. The user feedback panel is used to set Facets and choices including the range of crawling sites, the minimum size of the word cloud and the word frequency rating. The scope of crawling pages is the list of domains on which the content of papers can be identified. The minimum frequency of the word cloud corresponds to the average frequency of the terms shown in the word cloud, one of the outcomes of information analysis. Note that the word "cloud" is an image of various words, each of them of different sizes. The regular word classification is a frequency metric for evaluating the terms in the word frequency map. After the data analysis, the result view panel is used to provide data view results. In the predefined knowledge panel, words are added or removed from the DB knowledge DB. The portion of the Web Cruiser involves data collector and file creator parts. The information collector practices sorting, aid for vocabulary and collection. Parsing includes scanning papers' urls or browsing articles' text information on certain blogs. For this function, the webpage of items based on the user-defined input Facet are checked with a standard Resource Locator (Rel) to count the number of articles' web pages using their corresponding URL.

In fact, this generates a URL to crawl documents on the basis of the number of web pages for publications scanned through parsing. It then allows the relevant URL to be used as a hypertext markup (HTML) document for all posts in the database. Language aid requires UTF-8 encoding to avoid internet crawl data loss. Extracts text information from the HTML archive from the posts. The author of the database generates a derived data file and stores it in the DB component's ripple Database. There are three feature blocks in the data analyzer component: preprocessor, analyzer and visualize. The preprocessor extracts the collected text details of the papers contained in the crawling Database beforehand. By comparing the text data of the articles to a set of words in predefined knowledge, the preprocessor classifies these papers into relevant, non-relevant articles. If one of the terms is included in the document, it is defined as a corresponding article by the preprocessor feature section. In the DB research element review, the text information of the relevant articles are saved. In addition, non-relevant papers are discarded with text information.

Two steps are taken: selection of terms and sorting of phrases. In the extraction stage vocabulary is derived from a text information of the relevant articles, which consists of three or more characters. The unnecessary words (for example, 'A,' 'A,' 'The,' and 'The,') from the text data of the relevant articles are removed in the word filters step. The findings of the data analysis were saved in the DB element analyzer. The viewfinder visualizes information on the basis of the effects of data analysis and the alternative values

defined through the user input table. The output is visualized as word clouds and frequency curves in the data visualization and stored in the analyzer DB. The DB modules are creeping DB, predefined DB and DB review. The crawling data collected from the Internet by the internet crawlers is stored by the Crawling DB. A set of terms is placed by the client in the predefined information DB via the predefined knowledge table. DB research holds the text information for the relevant posts, results of data analytics and results of data show.

III. IMPLEMENTATION

The Facet configurations of the output and the application to scan the articles' webpage is performed in the team configuration box of the GUI. The field Query is used to pick the entry Facet in the category box and the button "Confirm" is used to check the articles' WebPages based on the Facet data. This application was provided for the inspection for papers' web pages. First of all, the Xml document with details is demanded about the total number of papers' web pages.

In order to obtain the maximum number of web pages of the posts, the text information whose allocation is < select category="search-header." The text data is then extract from the HTML document by xpathSApply). Last blanks and quotations are excluded from text information and text data is shown in the "Total Post" tab in a crawling community box (e.g., the maximum number of articles' web pages).

The scope of crawling pages is specified in the "Crawling Page" field in order to collect text information on posts. The internet crawling application is then performed with the "Crawling" key. The HTML file containing the content of posts on a given website and the document information with a div="story Body" p > attribute meaning will be retrieved by the XPathS Apply) (feature from the HTML file as you crawl through the site.

The whole list of crawling pages set by the client is constantly carried out with the site crawling method. The text information for all the papers published were eventually collected and stored in the DB. In the "Review" button, the information preprocessing and data analysis is done. The predefined information is used for identification of papers in software preprocessing.

The terms can be inserted and omitted by the client in predefined information. In the predefined vocabulary community container, the "Insert" or "Delete" buttons are used to insert or delete terms from predefined information respectively. The str detect) (function tests whether the terms are found in the predefined information in the collected articles text content. In the case of a word contained in the article, the relevant article shall be classified as an article, or else it shall be classified as an article of no relevance. DB is used in research to store text information for relevant articles, thus non-relevant papers are excluded from the text data. Once data analysis is finished, information analytics are conducted to retrieve and filter terms.

For word extraction, the supply (method is used to extract terms consisting of three or more characters from the text information of the relevant articles. For word processing, the

text information of relevant articles with the screen detect) (feature delete unnecessary words. The data analysis findings are contained in the DB archive. Visualizing the information is done using the "View" key. In the form of words or term intensity charts, the effects of the data visualization are shown. This is done using the choice values in the team option box (i.e., term intensity rating and minimum word cloud size). Word clouds are created using the word cloud () (function and the bar plot) function creates text frequency graphs. The DB research holds the output of data visualization.

IV. RESULT ANALYSIS

A word frequency chart showing 10 of the most frequently extracted words from the articles with 10 words. The graph's x and y axes are the word forms and terms rate, respectively. The frequency and intensity values of every term show on the chart on the line. In the figure, 'cloud' frequency is 2132, and 'cloud' frequency is 2%.

Parameter Values	Facet Data Analysis
Crawling page	300
Predefined	knowledge Data, IoT, Hadoop, Cloud
Filtering words	Good, will, then, the, are, with, and, that, this, but, have, has, can, for, you, from, been, more, they, said, what, its, about, how, was, which, their, into, these, when, there, those
Word cloud minimum frequency	150
Word frequency ranking	10

Table 1 Experimental parameters.

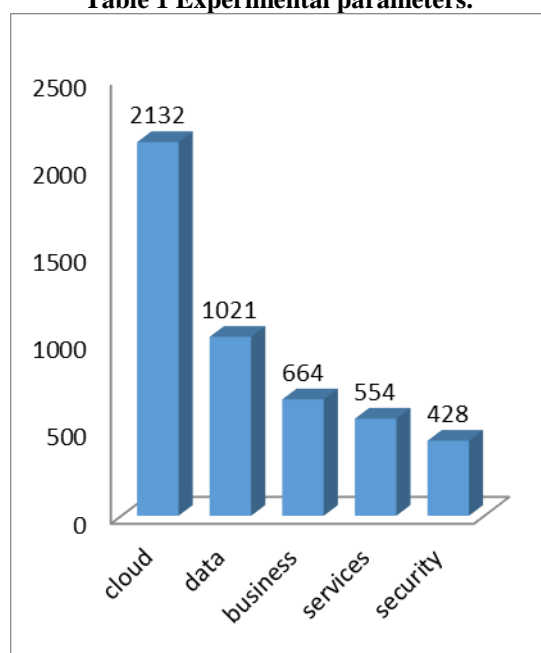


Fig1: Word frequency graph

V. CONCLUSIONS

In this paper they presented an automated site crawling and computational model for the processing and evaluation of Facet-based text content. Four components comprise the integrated frameworks: user interface, web crawling, data analyzation and DB. The user interface component helps the user to set the input Facet and the GUI option values. And the input Facet text data is compiled by a webcasting component. The information analyzer is used to preprocess, analyze data and model data. Finally, the DB component saves the text data collected, knowledge predefined, analysis results and results of data visualization. PoC prototyping was conducted to check the feasibility of the integrated framework. For the experiment, the text data of ZDNet items have been collected using the user's input Facet and text data gathered to monitor their frequencies have been analyzed. In a words-cloud and word frequency chart, the analysis results were visualized. The results showed that the integrated framework provides web crawling and text data analysis with reliable functions. The results show.

References

- [1] C. Dobre and F. Khafa: *Future Gener. Comp. Syst.* 37 (2014) 267.
- [2] W. Raghupathi and V. Raghupathi: *Health Inf. Sci. Syst.* 2 (2014).
- [3] Z. Khan, A. Anjum, and S. L. Kiani: *Proc. 2013 IEEE/ACM sixth Int. Conf. Utility and Cloud Computing (IEEE, 2013)* 381.
- [4] A. Sheth, C. Henson, and S. S. Sahoo: *IEEE Internet Comput.* 12 (2008) 78.
- [5] S. Landset, T. M. Khoshgoftaar, A. N. Richter, and T. Hasanin: *J. Enormous Data* 2 (2015) 24.
- [6] F. Morandat, B. Slope, L. Osvald, and J. Vitek: *Proc. European Conf. Item Oriented Programming (Springer, 2012)* 104.
- [7] B. C. D. S. Oliveira and J. Gibbons: *J. Funct. Program.* 20 (2010) 303.
- [8] A. Mesbah, A. V. Deursen, and S. Lenselink: *ACM Trans. Web* 6 (2012) 1.
- [9] Y. Zhang: *IEEE Trans. Serv. Comput.* 9 (2016) 786.
- [10] S. Wang, C. Zhang, and D. Li: *Proc. Int. Conf. Modern IoT Technologies and Applications (Springer, 2016)*
- [11] R Foundation: <https://www.r-project.org/> (got to July 2017).
- [12] Oracle: <http://www.oracle.com/technetwork/java/javase/downloads/jdk8-downloads-2133151.html> (got to July 2017). ZDNet: <http://www.zdnet.com/> (got to July 2017)