_____

# Survey on Mining Effective Information Using Ontology Based Semantic Web Crawler Mechanism

Mr. Bandar Muneer Khan Aslam
Department of Computer Engineering
P.V.P.I.T. College, Bavdhan
e-mail-muneerkhan.er@gmail.com

Mr. Gurav Yogesh B.
HOD
Department of Computer Engineering
P.V.P.I.T College, Bavdhan
e-mail-ybgurav@gmail.com

**Abstract—** Due to usable of copious data on web, searching has a consequential impact. Ongoing study place emphasis on the relevancy and robustness of the data found, as the invent patterns proximity is far from the probe. In spite of their relevance pages for some investigate topic, the results are mammoth that needed and are explored. Also the users' perspective differs in timely manner from topic to topic. In general terms ones' want is others unnecessary. Crawling algorithms play crucial role in selecting the pages that satisfies the users' needs. This paper reviews the research work on web crawling algorithms used on searching.

_____**\*\*\*\*\***_____

## I. INTRODUCTION

Text mining entirely utilized to achieving unknown information from natural language operationing and data supplying by applying various techniques.

In this technique for discovering the importance of term in document, term frequency of term is calculated. Sometime we can poster that two terms having same frequency in document but one term precede more meaning than other, for this concept based mining model is intended. In proffer model three measures are evaluated for analyzing concept in sentence, document and corpus levels. Semantic role labeler is mostly associated with semantic terms. The term which has more semantic role in sentence, it's known as Concept. And that Concept may be either word or phrase depending on sentence of semantic structure. When we put new document in system, the proposed model discover concept match by scanning all new documents and take out matching concept. Similarity measure used for concept analysis on sentence, document and corpus level that exceeds similarity measures depending on the term analysis model of document. The results are measured by using F-measure and Entropy. This model we are going to use for web documents.

## II. LITERATURE SURVEY

In this part, we briefly introduce the fields of semantic focused crawling and ontology-learning-based cored crawling, and review previous work on ontology learning-based cored crawling. A semantic focused crawler is it a software agent that is capable to traverse the Web, and retrieve correctly download related Web information on specific on specific topics by intend of semantic technologies [11], [12][22]. Since semantic way technologies provide shared knowledge for enhancing the interoperability among heterogeneous components, semantic technologies have on been broadly applied in the field of industrial automation [13]–[15]. The goal of semantic focused crawlers is to precisely and efficiently retrieve and download relevant Web information by automatically understanding the semantics underlying the Web information and the semantics underlying the predefined topics.

A survey actiond by Dong *et al.* [16] found that most of the crawlers in this domain make use of ontologies to represent the knowledge underlying topics and Web corroborate. However, the limitation of the ontology-based semantic focused crawlers is that the crawling performance essential depends on the attribute of ontologies. Further on, the features of ontologies may be altered by two issues. The first edition states, as it is well known that an ontology is the formal representation of specific domain knowledge [17] and ontologies are designed by domain master, a discrepancy may exist between the domain experts' understanding of the domain knowledge and the domain knowledge that exists in the real world. The second edition is that knowledge is dynamic and is changeable evolving, compared with relatively static ontologies. These two paradoxical site could lead to the enigma that ontologies sometimes cannot precisely represent real-world knowledge, pondering the issues of differentiation along with dynamism. The reproach of this problem in the field of semantic focused crawling is that the ontologies used by semantic focused crawlers cannot precisely represent the knowledge revealed in Web information, because Web information is mostly created or updated by human users with different knowledge understandings, and human intellectuals are effectual learners of new knowledge. The eventual issue of this problem is reflected in the gradually descending curves in the performance of semantic focused crawlers.

_____

In order to resolve the defects in ontologies and continue or enhance the performance of semantic-focused crawlers, study have begun to pay attention to enhancing semantic- focused crawling technologies by integrating them with ontology learning technologies. The goal of ontology erudition is to semi-automatically extract facts or patterns from a corpus of data and turn these into machine-readable ontologies [18].

Various techniques have been designed for ontology erudition, such as statistics-based techniques, linguistics (or natural language operation)-based techniques, logic-based techniques, etc. These techniques can also be categorize into manage techniques, semi-supervised techniques, and unsupervised techniques from the perspective of learning control. Plain, ontology-learning-based techniques can be used to solve the issue of semantic-focused crawling, by erudition new knowledge from crawled documents and merge the new knowledge with ontologies in order to constantly refine the ontologies.

## 2.1 WEB CRAWLER STRATEGIES:

### 2.1.1 Breadth First Search Algorithm:

This algorithm aims in the uniform investigate across the neighbour nodes. It institute at the root node and searches the all the neighbour nodes at the different level. If the objective is reached, then it is reported as success and the search is conclude. If it is not, it proceeds down to the next level sweeping the computation time. bEMADS and gEMADS these two algorithms are used based on Gaussian mixture model. Both resumes data into sub cluster and after that generate Gaussian mixture. These two algorithms run several orders of magnitude faster than maximum with little loss of search across the neighbour nodes at that level and so on until the object is reached. When all the nodes are investigate, but the objective is not met then it is reported as failure.

Breadth first is well suited for situations where the objective is found on the shallower parts in a deeper tree. It will not perform so well when the offshoot are so many in a game tree, signally like chess game and also when all the path leads to the same objective with the same length of the path [7][8].

Andy yoo et al [9] proposed a distributed BFS for numerous branches using Poisson random graphs and achieved high scalability through a set of clever memory and communication optimizations.

### 2.1.2 Depth First Search Algorithm

This powerful technique of systematically traverse through the search by starting at the root node and cross deeper through the child node. If there are multifarious than one child, then priority is offer to the left most child and traverse deep until no more child exists that seems to be available. It is backtracked to the next unvisited node and then stand in a similar manner [10].

This algorithm makes sure that all the edges are visited once breadth [11]. It is well suited for problems regarding searches, but when the branches extends in length then this algorithm takes might end up in an infinite loop [8].

### 2.1.3 Page Rank Algorithm

Page rank algorithm determines the importance of the web pages by counting citations or back links to a given page [12]. The page rank of a given page is calculated as

$$PR(A) = (1-d) + d (PR(T1)/C(T1) + ... + PR(Tn)/C(Tn))$$

PR(A) □Page Rank of a Website,

d □damping factor

$T_1,....T_n$ □links

Yongbin Qin and Daoyun Xu [13] proffer an algorithm, taking the human factor into reflection, to introduce page belief recommendation mechanism and brought forward a balanced rank algorithm based on PageRank and Page belief recommendation which ultimately attaches importance into the subjective needs of the users; so that it can effectively avoid topic drift problems. Tian Chong [14] proposed a new type of algorithm of page distinction by compound classified tree with static algorithm of PageRank, which enables the categorize tree to be form according to a large number of users' similar searching results, and can plain reduce the problem of Theme-Drift, caused by using PageRank only, and puzzle of outdated web pages and extend the efficiency and skillful of search. J.Kleinberg [15] proposed a dynamic page ranking algorithm. Shaojie Qiao [16] proposed a new page rank algorithm based on similarity measure from the vector space paragon, called SimRank, to score web pages. They proffer a new similarity measure to compute the similarity of pages and apply it to partition a web database into several web social networks (WSNs)

### 2.1.4 Genetic Algorithm

Genetic algorithm is basis on biological progression whereby the fittest offspring is obtained by crossing over of the selection of some best individuals in the population by intend of fitness function. Therefore a search algorithm solutions to the problem exists but the technique is to find the best solution within specified time [17]. [18] Shows the genetic algorithm is best suited when the user has literally no or less time to spend in searching a huge database and also very efficient in multimedia results. While almost all conventional methods search from a single point, Genetic Algorithms always operates on a whole population. This provide much to the robustness of genetic algorithms. It eventually reduces the risk of becoming trapped in a local stationary point [19].Now

further on the fitting of Genetic Algorithms by various researchers [23], [24], [25], [26] has been reproduce in [27].

### 2.1.5 Naïve Bayes classification Algorithm

Naïve Bayes algorithm is based on Probabilistic learning and classification. It follows the assumption that one feature is independent of another [28]. This algorithm proved to be efficient over many other approaches [29] although its simple assumption is not much applicable in realistic cases [28]. Wenxian Wang et al [30] proposed an efficient crawler based on Naïve Bayes to gather many relevant pages for hierarchical website layouts. Peter Flach and Nicolas Lachiche [31] presented Naïve Bayes classification of structured data on artificially generated data.

### 2.1.6 HITS Algorithm

This algorithm put forward by Kleinberg is previous to Page rank algorithms which uses scores to calculate the relevance [32]. This method retrieves a set of results for a search and calculate the authority and hub score within that set of results. Because of these reasons this method is not often used [2].

### III. CONCEPTS HELPFUL TO MEASUREMENT OF SIMILARITIES

The techniques prescribed in the previous work are used for document clustering. But it is only for documents present on system. In the proposed system we are going to use web documents and we will get the clustered output and that have shown in fig.1. This concept-based mining model consists of sentence-based concept analysis, document-based concept analysis, corpus-based concept-analysis, and concept-based similarity measure.. A web document is given as the input to the proposed model. Each document has well-defined sentence boundaries.  Each sentence in the document is labeled automatically based on the Prop Bank notations. After running the semantic role labeler, labeled verb argument structures, the output of the role labeling task, are captured and analyzed by the concept-based mining model on sentence, document, and corpus level. In the concept-based mining model, a labeled terms either word or phrase is considered as concept.

The proposed model contains the following modules

#### A. Web Document

Web document is given as Input to the given system. Here user can give any query to the browser. Pure HTML pages are selected by removing extra scripting. Web pages contain data such as hyperlinks, images, script. So it is necessary to remove such unwanted script if any, during the time when a page is selected for processing. The HTML code is then transferred into XML code. On that document next process is processed that is Text pre-processing or data processing.

#### B. Data Processing

First step is separate sentences from the documents. After this label the terms with the help of Prop Bank Notation. With the help of Porter algorithm remove the stem word and stop words from the terms.

#### C. Concept Based Analysis

This is important module of the proposed system. Here we have to calculate the frequencies of the terms.

Conceptual term regularity (ctf), Term frequency (tf) and Document regularity (df) are calculated. The objective behind the concept-based analysis task is to achieve an true analysis of concepts on the sentence, document, and corpus levels rather than document only.

##### 1) Sentence based concept analysis

For analyzing every concept at sentence level, concept based frequency measure; called conceptual term frequency is used.

##### 1.1) Calculating ctf in sentence s

Ctf is the number of occurrences of concept c in verb structure of sentence s. If concept c frequently appears in structure of sentence s then it has principal role of s.

##### 1.2) Calculating ctf in document d

A concept c can have many ctf values in different sentences in the same document d. Thus, the ctf value of concept c in document d is calculated by:

$$ctf = \frac{\sum_{n=1}^{sn} ctf_n}{sn},$$

Where sn: total number of sentences containing concepts in document d

##### 2) Document Based Concept Analysis

For analyzing concepts at document level term frequency tf in original document is calculated. The tf is a local measure on the document level.

##### 3) Corpus Based Concept Analysis

To calculate concepts from documents, document frequency df is used. Document Frequency df is the global measure. With the help of Concept based Analysis Algorithm we can calculate ctf, tf, df.

#### D. Similarity Approach

This module mainly contains three parts. Concept based similarity, Singular Value Decomposition and combined based similarity it contains. Here we get that how many percentage of concept math with the given web document.

*E. Concept Based Similarity*

A concept-based similitude measure depends on matching concept at sentence, document, and corpus instead of individual terms. This similitude measure based on three main appearance. First is explicate label terms that capture semantic structure of each sentence. Second is concept frequency that is used to measure participation of concept in sentence as well as document. Last is the concepts measured from number of documents.

Concept based similarity between two document is calculated by:

$$sim_c(d_1, d_2) = \sum_{i=1}^{m} max\left(\frac{l_{i_1}}{Lv_{i_1}}, \frac{l_{i_2}}{Lv_{i_2}}\right) \times weight_{i_1} \times weight_{i_2}$$

Term frequency is calculated by following formula:

$$tf\ weight_i = \frac{tf_{ij}}{\sqrt{\sum_{j=1}^{cn}(tf_{ij})^2}},$$

*F. Clustering Techniques*

This module used three main basic techniques like Single pass, Hierarchical Agglomerative Clustering, and K-Nearest Neighbor. With the help of these techniques we can get that which cluster is having highest priority.

*G. Output Cluster*

Last module is the output Cluster. After applying the clustering techniques we get clustered document. That will help to find out main concepts from the web document.

## IV. SYSTEM IMPLEMENTATION

The proposed system model illustrates flow of implementation. First, web document given input to the system where, HTML pages are collected and their XML conversion is carried out. In the second module that is in Text Processing carried out separate sentences, label terms, and removing stop words and stem words. Third module Concept based analysis measures conceptual term frequency (ctf), term frequency (tf), and document frequency (df). Next module concept based document similarity find out how many percentage of concept is similar to the given concept.

## V. CONCLUSION AND FUTURE WORK

The main objective of the review paper was to throw some light on the web crawling algorithms. We also discussed the various search algorithms and the researches related to respective algorithms and their strengths and weaknesses associated. We believe that all of the algorithms surveyed in this paper are effective for web search, but the advantages

favours more for Genetic Algorithm due to its iterative selection from the population to produce relevant results

## VI. REFERENCES

[1] Pavalam S M, Jawahar M, Felix K Akorli, S V Kashmir Raja " Web Crawler in Mobile Systems" International Conference on Machine Learning (ICMLC 2011), Vol. , pp

[2] Alessio Signorini, "A Survey of Ranking Algorithms" retrieved from http://www.divms.uiowa.edu/~asignori/phd/report/asurvey-of-ranking-algorithms.pdf 29/9/2011

[3] Maurice de kunder, "Size of the world wide web", retrieved from http://www.worldwidewebsize.com/ 8/8/11

[4] Ricardo Baeza-Yates, Ricardo Baeza-Yates "Crawling a Country: Better Strategies than Breadth-irst for Web Page Ordering" , Proc. WWW 2005.

[5] Marc Najork, "Web Crawler Architecture" retrieved from http://research.microsoft.com/pubs/102936/EDSWeb CrawlerArchitecture. pdf accessed on 10/8/11

[6] Junghoo Cho and Hector Garcia-Molina "Effective Page Refresh Policies for Web Crawlers" ACM Transactions on Database Systems, 2003.

[7] Steven S. Skiena "The Algorithm design Manual" Second Edition, Springer Verlag London Limited, 2008, Pg 162.

[8] Ben Coppin "Artificial Intelligence illuminated" Jones and Barlett Publishers, 2004, Pg 77.

[9] Andy Yoo, Edmond Chow, Keith Henderson, William McLendon, Bruce Hendrickson, ÄUmit CatalyÄurek "A Scalable Distributed Parallel Breadth-First Search Algorithm on BlueGene/L" ACM 2005.

[10] Alexander Shen "Algorithms and Programming: Problems and solutions" Second edition Springer 2010, Pg 135

[11] Narasingh Deo "Graph theory with applications to engineering and computer science" PHI, 2004 Pg 301

[12] Sergey Brin and Lawrence Page "Anatomy of a Large scale Hypertextual Web Search Engine" Proc. WWW conference 2004

[13] Yongbin Qin and Daoyun Xu "A Balanced Rank Algorithm Based on PageRank and Page Belief recommendation"

[14] TIAN Chong "A Kind of Algorithm For Page Ranking Based on Classified Tree In Search Engine" Proc International Conference on Computer Application and System Modeling (ICCASM 2010)

[15] J.Kleinberg "Authoritative sources in a hyperlinked environment", Proc 9th ACM-SIAM Symposium on Discrete Algorithms, 1998.

[16] Shaojie Qiao, Tianrui Li, Hong Li and Yan Zhu, Jing Peng, Jiangtao Qiu "SimRank: A Page Rank Approach based on similarity measure" 2010 IEEE

[17] S. N. Sivanandam, S. N. Deepa "Introduction to Genetic Algorithms" Springer, 2008, pg 20

[18] S.N. Palod, Dr S.K.Shrivastav, Dr P.K.Purohit "Review of Genetic Algorithm based face recognition" International Journal of Engineering Science and Technology (IJEST) Vol. 3 No. 2 Feb 2011

[19] Deep Malya Mukhopadhyay, Maricel O. Balitanas, Alisherov Farkhod A., Seung-Hwan Jeon, and Debnath

_____

Bhattacharyya "Genetic Algorithm: A Tutorial Review" International Journal of of Grid and Distributed Computing Vol.2, No.3, September, 2009.

[20] Shian-Hua Lin, Jan-Ming Ho, Yueh-Ming Huang , ACRID , intelligent internet document organization and retrieval , IEEE Transactions on Knowledge and data engineering, 14(3), 559-613, 2002.

[21] Pavalam S M1, S V Kashmir Raja2, Felix K Akorli3 and Jawahar M4, " A Survey of Web Crawler Algorithms" National University of Rwanda Huye, RWANDA

[22] Hai Dong, Member, IEEE, and Farookh Khadeer Hussain" Self-Adaptive Semantic Focused Crawler for Mining Services Information Discovery" IEEE TRANSACTIONS ON INDUSTRIAL INFORMATICS, VOL. 10, NO. 2, MAY 2014.

[23] Zacharis Z. Nick and Panayiotopoulos Themis, Web Search Using a Genetic Algorithm, IEEE Internet computing, 1089-7801/01c2001, 18-25, IEEE

[24] Ramakrishna Varadarajan, Vagelis Hristidis, and Tao Li , Beyond Single-PageWeb Search Results, IEEE Transactions on knowledge and data engineering, 20(3) , 411 - 424, 2008

[25] Judit BarIlan, Comparing rankings of search results on the Web, Information Processing and Management 41 (2005) 1511–1519

[26] Adriano Veloso, Humberto M. Almeida, Marcos Goncalves, Wagner Meira Jr., Learning to Rank at Query-Time using Association Rules, SIGIR'08, 267-273 , 2008, Singapore.

[27] S.Siva Sathya and Philomina Simon, " Review on Applicability of Genetic Algorithm to Web Search" International Journal of Computer Theory and Engineering, Vol. 1, No. 4, October2009

_____