

# An Effective Web Ontology Using Web Crawler Systems to Measures Web Similarities

M.Florence Dayana

Ph.D., Research Scholar,

A.V.V.M Sri Pushpam College (Autonomous), Poondi, Thanjavur, Tamilnadu.

Email: florencedayana@gmail.com

Dr.M.Chidambaram

Assistant Professor in Computer Science,

Rajah Serfoji Government College (Autonomous), Thanjavur, Tamilnadu.

Email: chidsuba@gmail.com

**Abstract**— Information mining is the process of extraction of hidden predictive information from the colossal databases. It is a new innovation with great latent to help companies focus on the most essential information in their information warehouses. Web mining is an information mining strategies which naturally find information from web documents. A focused crawler may be portrayed as a crawler which returns pertinent web pages on a given point in navigating the web. Web Crawlers are one of the most crucial part of the Seek Engines to collect pages from the Web. Service-Oriented Computing (SOC) is an interdisciplinary paradigm that revolutionizes the exceptionally fabric of appropriated programming advancement applications that adopt Service-Oriented Architectures (SOA) can evolve amid their lifespan and adapt to changing or unpredictable environments more easily. The framework incorporates the advancements of semantic focused crawling and ontology learning, in demand to maintain the execution of this crawler. The crawler we outlined is used to solve the issue relating to the heterogeneity, ubiquity and ambiguity and machine learning will increment the execution of the crawler and moreover perform prediction on data. This paper criticisms investigates on web crawling procedures used on searching web similarities.

**Keyword** - Ontology Learning; Semantic Focused Crawler; Web Mining; Similarity Measures; Web Crawler;

## I. INTRODUCTION

A crawler is an operator which can naturally seek and download WebPages. Focused (topical) crawlers are a bunch of appropriated crawlers that practice in certain particular topics. Each crawler will dissect its topical limit when bringing WebPages.

A semantic focused crawler is a programming operator that is capable to navigate the Web, and recover as well as download related Web information for particular topics, by implies of semantic Web technologies. The objective of semantic focused crawlers is to accurately and effectively recover and download pertinent Web information by understanding the semantics fundamental the Web information and the semantics fundamental the prestructured topics.

After the World Wide Web emerged, scientists endeavored to improve its quality by different semantic technologies. Currently there are three new frames of perceived networks improved by different semantic technologies, which are semantic web, semantic grid, and information grid. Semantic web is “a web of data”, which is used to express the meaning of web information by implies of various ontological mark-up languages, such as XML, RDF, OWL and so forth. It gives the machine-capable information for PCs to retrieve, offer and blend information on the internet.

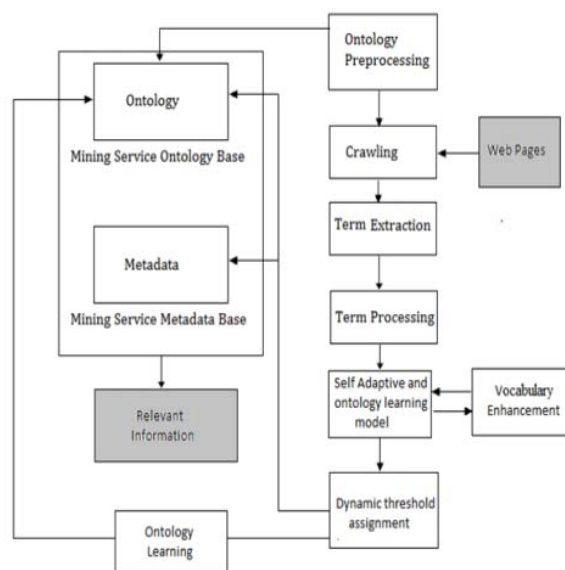


Fig.1 Web Mining Ontology

### A. Overview of Web Mining

Web mining alludes to the discovery of information from Web information that incorporate Web pages, media objects on the Web, Web links, Web log data, and other information generated by the use of Web data. Web mining is classified into: (a) Web content mining, (b) Web structure mining and (c) Web use mining. Web content mining alludes to mining information from Web pages and other Web objects. Web structure mining alludes to mining information about link structure connecting Web pages and other Web objects. Web use mining alludes to the mining of use patterns of web pages found among clients accessing a Website. Among the three, Web content mining is perhaps studied most extensively due to the prior work in content mining. The traditional topics covered by Web content mining include:

#### i) Web page arrangement

This includes the arrangement of Web pages under some pre-defined categories that may be organized in a tree or other structures.

#### ii) Web clustering

This includes the grouping of Web pages based on the similarities among them. Each resultant group should have similar Web pages while Web pages from distinctive resultant groups should be dissimilar.

#### iii) Web extraction

This includes extracting HTML elements, term phrases, or tuples from Web pages that represent some required idea instances. The popularity of WWW is largely dependent on the look engines. Look motors are the gateways to the huge information repository at the internet. Look engine consist of four discrete components: Crawling, Indexing, Ranking and query-processing. The earliest Web look motors relied for retrieving information from Web pages on the bases of coordinating with the words in the look query. As the Web continues to grow, and as the diversity of clients increases look motors must utilize semantic clues to satisfy user's information needs. Semantic search, which takes into account the interests of the client as well as the specific concontent in which the look is issued, is the next step in providing clients with the most important information possible. Currently the general purpose look motors strive as entry points for the web pages perform the coverage of information that is as broad as possible. They use Web crawlers to maintain their index databases. These crawlers are blind and exhaustive in their approach, with comprehensiveness as their major goal. A URL (Uniform Resource Locator) that is a URI (Uniform Resource Identifier) specifies where an identified Asset is available. In order to look most important information, crawlers can be more selective about the URL they fetch and refer as to be crawled this mechanism.

The Web has becoming the biggest unstructured database for accessing information over the documents. It is well recognized that the information technology has a profound effect on the conduct of the business, and the Web has become the biggest marketplace in the world. Innovative business professionals have realized the commercial applications of the Web for their customers and strategic partners. With the rapid growth of electronic content from the complex the WWW, more and more information you need is included. But, the massive sum of content also takes so much trouble to people to find useful information. For example, the standard Web look motors have low precision, since typically some important Web pages are returned mixed

with a large number of important pages, which is fundamentally due to the situation that the topic-specific features may occur in distinctive contexts. So, one appropriate way of organizing this overwhelming sum of reports is necessary. The World Wide Web is an architectural framework for accessing linked reports spread out over millions of machines all over the Internet.

## II. RELATED WORK

### A. Semantic Focused Web Crawler

A crawler is an operator which can naturally seek and download WebPages. Focused (topical) crawlers are a bunch of appropriated crawlers that practice in certain particular topics. Each crawler will dissect its topical limit when bringing WebPages. A semantic focused crawler is a programming Operator that is capable to navigate the Web, and recover as well as download related Web information for particular topics, by implies of semantic Web technologies. The objective of semantic focused crawlers is to accurately and effectively recover and download pertinent Web information by understanding the semantics fundamental the Web information and the semantics fundamental the precharacterized topics.

After the World Wide Web emerged, scientists endeavored to improve its quality by different semantic technologies. Currently there are three new frames of perceived networks improved by different semantic technologies, which are semantic web, semantic grid, and information grid. Semantic web is “a web of data”, which is used to express the meaning of web information by implies of various ontological markup languages, such as XML, RDF, OWL and so forth. It gives the machine capable information for PCs to retrieve, offer and blend information on the internet.

### B. Specialists and Web Administrations

The use of programming specialists in Web Administration disclosure has been the subject of numerous researches. Creators in proposed a multi-Operator approach to achieve administration disclosure and determination agreeing to the consumers' preferences. They portray a framework in which specialists interact and offer information, in essence creating an Eco framework of collaborative administration suppliers and consumers. In, creators propose a multi-Operator approach for a appropriated information recovery task. In this work, each Operator has a view of its environment called Operator view.

The agent-view structure of an Operator contains information about the language models of records owned by each agent. An agent-view reorganization algorithm is run to dynamically reorganize the fundamental agent-view topology. The proposed protocol does not use ontologies amid information retrieval.

## III. WEB CRAWLER STRATEGIES

### A. Broadness First Search

Broadness First Seek is the simplest form of crawling algorithm. It starts with a join and keeps on navigating the connected joins without taking into consideration any information about the topic. Since it does not take into account the relevancy of the path while traversing, it is moreover known as the Blind Seek Algorithm. It is considered to give lower bound on effectiveness for any intelligent traversal algorithm.

### B. Page Rank Algorithm

PageRank is the connectivity-based page quality measure. It is a static measure; it is outlined to rank pages in the absence of any queries. That is, PageRank computes the “global worth” of each page. Intuitively, the PageRank measure of a page is comparable to its in-degree, which is a possible measure of the importance of a page. The PageRank of a page is high if numerous pages with a high PageRank contain joins to it, and a page containing few active joins contributes more weight to the pages it joins to than a page containing numerous active links. The page rank of a given page is calculated as

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

PR(A) Page Rank of a Website,

d damping factor

T1,...,Tn links

### C. Hereditary Algorithm

Hereditary algorithm is based on biological evolution whereby the fittest offspring is obtained by crossing over of the determination of some best individuals in the populace by implies of fitness function. In a seek algorithm solutions to the issue exists but the technique is to find the best solution within specified time Shows the Hereditary algorithm is best suited when the client has literally no or less time to spend in searching a colossal database and moreover exceptionally effective in multimedia results. While almost all conventional strategies seek from a single point, Hereditary Calculations always operates on a entirety population. This

contributes much to the robustness of Hereditary algorithms. It reduces the risk of becoming trapped in a nearby stationary point.

#### D. Naive Bayes Algorithm

Naive Bayes algorithm is based on Probabilistic learning and classification. It assumes that one feature is independent of another. This algorithm proved to be effective over numerous other approaches although its straightforward assumption is not much applicable in realistic cases. An effective crawler based on Naive Bayes to gather numerous pertinent pages for hierarchical website layouts. Naive Bayes arrangement of structured information on artificially created data.

#### E. HITS Algorithm

Hyperlink-Induced Point Seek is a join analysis algorithm that rates Web pages, developed by Jon Kleinberg. The idea behind Hubs and Authorities stemmed from a particular insight into the creation of web pages when the Web was originally forming; that is, certain web pages, known as hubs, served as large directories that were not actually definitive in the information that they held, but were used as compilations of a broad catalog of information that led users direct to other definitive pages. In other words, a great hub reintroduced a page that pointed to numerous other pages, and a great authority reintroduced a page that was connected by numerous different hubs.

### IV. PROPOSED APPROACHES

The future scheme typical demonstrates movement of execution. Mainly, web file assumed input near the alteration is agreed out. In the next segment that is in Text Dispensation agreed out single sentences, label terms, plus eliminating break words and stalk arguments. Third segment perception created investigation actions theoretical term incidence (tti), term incidence (ti), and script incidence (si). Then the last segment built web text similarity discover out how various measurement of perception is similar toward the agreed model. Structure wherever, HTML sheets be located and collected their XML.

#### A. Framework Evaluation

Five execution indicators of information recovery to evaluate ontology based focused crawler are used which are harvest rate, precision, recall, consonant mean, aftermath rate and crawling time. Harvest rate in the information recovery is used to measure the crawling ability of a crawler. Harvest rate is the proportion of coherently connected metadata to the entirety gathering of metadata. It can be characterized below:

$$\text{Harvest rate} = \frac{\text{No. of logically linked MetaData}}{\text{No. of MetaData}}$$

It can be seen that the proposed model has the optimal execution (52%), compared to the semantic model (around 46%). Exactness in the information recovery is used to measure the preciseness of a recovery system. Exactness for a single idea is the proportion of coherently connected and pertinent metadata to all coherently connected metadata. It can be characterized below:

$$\text{Precision}(S) = \frac{\text{No. of logically linked relevant MetaData}}{\text{No. of logically linked MetaData}}$$

The entirety exactness is the whole of exactness for each idea in the gathering and it can be characterized below:

$$\text{Precision}(W) = \frac{\sum_{i=1}^n \text{precision}(S_i)}{n}$$

It can be observed that the overall exactness of the proposed model is 34.20%, and the overall exactness of the semantic model is 31.59%. Review in the information recovery is used to measure the effectiveness of a query system. Review is the proportion of coherently connected and pertinent metadata to all pertinent metadata. It can be characterized below:

$$\text{Recall}(S) = \frac{\text{No. of logically linked relevant MetaData}}{\text{No. of logically linked MetaData}}$$

The absolute review is the whole of review for each idea in the collection. It can be characterized below:

$$\text{Recall}(W) = \frac{\sum_{i=1}^n \text{Recall}(S_i)}{n}$$

It can be seen that the overall review for the proposed model is 65.59%, compared to only 63.21% for the semantic model.

The aftermath rate is proportion of coherently connected and non-pertinent metadata to the entirety gathering of non-pertinent metadata to the concept. It is characterized below:

$$\text{Fallout rate}(S) = \frac{\text{No. of logically linked non - relevant MetaData}}{\text{No. of non - relevant MetaData}}$$

The entirety aftermath rate is the whole of aftermath rate for each idea in the collection. It can be reintroduced as:

$$\text{Fallout rate}(W) = \frac{\sum_{i=1}^n \text{Fallout rate}(S_i)}{n}$$

The fall rate esteem is lower, the crawler's execution is better. It can be seen that the overall aftermath rate of the proposed model is 0.43%, and for the semantic model is around 0.46%. Consonant Mean is characterized as the absolute execution of exactness and review value. Consonant Mean is define as below:

$$\text{Harmonic Mean} = \frac{\text{Precision} + \text{Recall}}{\text{Precision} \cdot \text{Recall}}$$

As an aggregated parameter, the overall consonant mean values for both of these models are beneath 50% (43% for the semantic model and 44.46% for the proposed model). Crawling time is used to measure the effectiveness of a crawler. The crawling time for a Web page is characterized as the time interval from processing the Web page from the Crawling process to the Metadata Generation and Association process or to the Separating process. Overall, the average crawling speed for the proposed model is 77.00 ms/page and for the semantic model is 70.00 ms/page.

#### B. System Implement to Measure the Similarities

The strategies prescribed in the previous work are utilized for report clustering. But it is only for reports present on system. In the proposed framework we are going to use web reports and we will get the grouped yield and that have shown in fig.1. This concept-based mining model consists of sentence-based idea analysis, document-based idea analysis, corpus- based concept-analysis, and concept-based closeness measure. A web report is given as the input to the proposed model. Each report has well-defined sentence boundaries. Each sentence in the report is named automatically based on the Prop Bank notations. After running the semantic part labeler, named verb argument structures, the yield of the part 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 named terms either word or phrase is considered as concept.

The proposed model contains the following modules

##### i) Web Report

Web report is given as Input to the given system. Here client can give any inquiry to the browser. Pure HTML pages are chosen by removing extra scripting. Web pages contain information such as hyperlinks, images, script. So it is essential to expel such unwanted script if any, during the time when a page is chosen for processing. The HTML code is then transferred into XML code. On that report next process is processed that is Content pre- handling or information processing.

##### ii) Information Handling

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

##### iii) Idea Based Investigation

This is important module of the proposed system. Here we have to figure the frequencies of the terms. Theoretical term frequency (tti), term frequency (ti), and script frequency (si) are calculated. The objective behind the concept-based investigation undertaking is to achieve an accurate investigation of ideas on the sentence, document, and corpus levels rather than report only.

###### a. Sentence based idea investigation

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

- Figuring tti in sentence s

tti is the number of occurrences of idea c in verb structure of sentence s. If idea c frequently appears in structure of sentence s then it has principal part of s.

- Figuring tti in report d

A idea c can have numerous tti values in distinctive sentences in the same report d. Thus, the tti value of idea c in report d is ascertained by:

*b. Report based idea investigation*

For analyzing ideas at report level term recurrence  $t_i$  in original report is calculated. The  $t_i$  is a local measure on the report level.

*c. Corpus based idea investigation*

To figure ideas from documents, report recurrence  $s_i$  is used. Report Recurrence  $t_i$  is the global measure. With the help of Idea based Investigation Algorithm we can figure tti,  $t_i$ ,  $s_i$ .

**iv) Closeness Approach**

This module fundamentally contains three parts. Idea based similarity, Singular Value Decomposition and combined based closeness it contains. Here we get that how numerous percentage of idea math with the given web document.

**v) Idea Based Closeness**

A concept-based closeness measure depends on coordinating idea at sentence, document, and corpus instead of individual terms. This closeness measure is based on three primary aspects. First are analyzed name terms that capture semantic structure of each sentence. Second is idea recurrence that is utilized to measure participation of idea in sentence as well as document. Last is the ideas measured from number of documents. Idea based closeness among two report is ascertained by:

**vi) Grouping Strategies**

This module utilized three primary basic strategies like Single pass, Hierarchical Agglomerative Clustering, and K- Nearest Neighbor. With the help of these strategies we can get that which Group is having highest priority.

**vii) Yield Group**

Last module is the yield Cluster. After applying the grouping strategies we get grouped document. That will help to find out primary ideas from the web document.

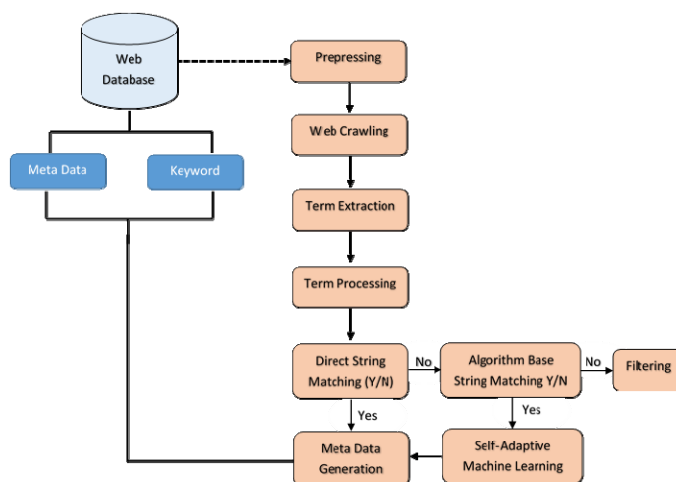


Fig.2 Similarity Measure Process

## V. CONCLUSION

This paper presents the overview of semantic focused web crawler and different strategies which are used for semantic focused web crawler for retrieving pertinent information from web. From above discussion it concluded that semantic focused web crawler has some limitations. So ontology learning based focused crawler is exceptionally valuable for finding pertinent information from web. We moreover discussed the different seek calculations and the researches related to respective calculations and their strengths and weaknesses associated. In this paper, we have design Self Modifying Semantic Focused Crawler to mined any kind of services. The foremost impartial investigation of the paper was to pitch certain bright on the mesh crawling procedures. We also conversed the several search algorithms and investigates connected to separate 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 (GA) due to its iterative selection from the population to produce relevant results.

## REFERENCES

- [1] Danushka Bollegala, Yutaka Matsu and Mitsuru Ishizuka, "A Web Search Engine-Based Approach to Measure Semantic Similarity among Words", in IEEE Transactions on Knowledge and Data Engineering, Vol- 23, Issue- 7, PP- 977 – 990, 2011.
- [2] Guosheng Kang, Mingdong Tang, Jianxun Liu, Xiaoqing (Frank) Liu and Buqing Cao, "Diversifying Web Service Recommendation Results via Exploring Service Usage History", in IEEE Transactions on Services Computing, Vol- 9, Issue-4, PP- 566 – 579, 2016.
- [3] Hai Dong and F. K. Hussain, "Self-Adaptive Semantic Focused Crawler for Mining Services Information Discovery", in IEEE Transactions on Industrial Informatics, Vol- 10, Issue-2, PP-1616 – 1626, 2014.
- [4] Huifeng Sun, Zibin Zheng, Junliang Chen and Michael R. Lyu, "Personalized Web Service Recommendation via Normal Recovery Collaborative Filtering", in IEEE Transactions on Services Computing, Vol-6, Issue-4, PP-573 – 579, 2013.
- [5] Hai Dong and F. K. Hussain, "Focused Crawling for Automatic Service Discovery, Annotation, and Classification in Industrial Digital Ecosystems", in IEEE Transactions on Industrial Electronics, Vol-58, Issue-6, PP-2106 – 2116, 2011.
- [6] Qingpeng Zhang and David Haglin, "Semantic similarity among ontologies at different scales", IEEE/CAA Journal of Automatica Sinica, Vol- 3, Issue- 2, PP- 132 – 140, 2016.
- [7] Sheau-Ling Hsieh, Wen-Yung Chang, Chi-Huang Chen and Yung-Ching Weng, "Semantic Similarity Measures in the Biomedical Domain by Leveraging a Web Search Engine", in IEEE Journal of Biomedical and Health Informatics, Vol- 17, Issue- 4, PP-853 – 861, 2013.
- [8] S. Wang, Q. Sun, H. Zou and F. Yang, "Reputation measure approach of web service for service selection", IET Software, Volume: 5, Issue: 5, PP-466 – 473, 2011.
- [9] Xuebo Song, Lin Li, Pradip K. Srimani and Philip S. Yu, James Z. Wang, "Measure the Semantic Similarity of GO Terms Using Aggregate Information Content", IEEE/ACM Transactions on Computational Biology and Bioinformatics, Vol- 11, Issue- 3, PP- 468 – 47, 2014.