# A Novel Approach for Semantic Based Deep Web Classification

M.Sreekrishna[1], P.Sundaramoorthy[2], Dr.T.Rajendran[3]

[1]*PG Scholar, CSE Department, Angel College of Engineering and Technology*
[2] *Assistant Professor, CSE Department, Angel College of Engineering and Technology*
[3]*Dean , CSE and IT Department, Angel College of Engineering and Technology*
[1]sree7krish@gmail.com, [2]sundarme08@gmail.com,[3]rajendran_tn@yahoo.com

*Abstract*— **The deep web are the web that are not a part of surface web. Due to the large volume of data deep web have grained a large attention in recent years. Recent search engines cannot be used to extract the information in the dark Web. The extension of deep web is found to be larger than the surface web. So classification of the deep web are necessary for to make the extraction relevant to the user. In this paper we propose a novel semantic ontology based deep web classification method SODWEB is used for to classify the data in the deep web automatically. The method is found to be highly efficient and scalable for textraction. The URL is chosen for the process of analyzing the semantic association among the concepts. Then that URL of deep web search source is mapped to the category hierarchy obtained. Generally, Wikipedia categories are taken for the process of analysis. We perform the process of finding the depth of the websites after the process of extraction and classification The experiments on the search is found to be accurate and fine grained on comparing to the results from keyword based automatic classification.**

*Keywords- Deep web, Information retrieval, Ontology, Semantic search, Wikipedia*.

## I. INTRODUCTION

World Wide Web is found to have surface web and deep web. Surface web are the web that are visible web whereas the deep web are the invisible web. The surface web consist of the portion of world wide web which is indexed by conventional search engines. The deep web is forund to be the content hidden behind the HTML forms. The information hidden in the deep web is found to more the surface web.

The estimated number [1] of deep web is found to be 43,000-96,000 with an estimated mass of 7,500 terabytes. The number shows that the volume of the deep web is 400 time larger than the data contained in the surface web.

Further the recent studies [2] have shown that the quality of deep web is found to better than the surface web. The relevancy factor is also high incase of the database on the deep web [3]. These database focus on the confined domain like health, sports, books, travel, hotels, property and vehicles etc. The content of these databases is naturally clustered which greatly improves user"s identification over searched content.

For to bring the hidden information sources to the surface web, the contents in these sources have to be crawled and indexed by the search engines[4-7].They were found to be not much efficient for the process of extracting the relevant web page. For to solve these problem research community proposes the automatic classification of hidden web sources into the domain categories as the first step[8-17].The articles in Wikipedia is not strictly categorized i.e. an article may belong to more than one categories. According to [18] there are almost 400,000 categories in the English Wikipedia, with an average of 19 articles and two subcategories each.

The key issue is to extract the relevant web pages in the web search. Existing web search is based on user"s query terms in which user issues the query topics keywords to search engine and it returns a set of pages that may be related to the query terms. It does not analyze the source code of the web page and only limited links can be retrieved. Recognizing these drawbacks, this work proposes a novel approach to the problem of retrieving relevant web pages.

The proposed system uses an efficient hyperlink based algorithm to find all relevant child links for a given web page. The input is found to be URL of a page and the output is found to be set of all child links.

## II. STRUCTURE OF DEEP WEB CLASSIFICATION

The two main approaches of deep web classification are manual and automatic. In manual approach classifications are done manually and the advantage is that it provides high quality of directory services. The disadvantage involves that human intervention which is found to be time consuming and expensive. Further it is found to be non-scalable.[19]

The database in the deep web are dynamic due to insertion and deletion mechanism so the database of the deep web have to updated automatically. So the automatic deep web classification was introduced.

Definition 1: Given a set of deep web data sources $D_i$ where i=1,…,n and a set of category domains $CD_i$ where i=1,…,n.A set of procedures $P_i$ where i=1,…,n are adopted to extract the content representative(CR) from each deep web sources. CR is

submitted to a classifier (C), which results in the classification of Di into one or more category domains CDi.
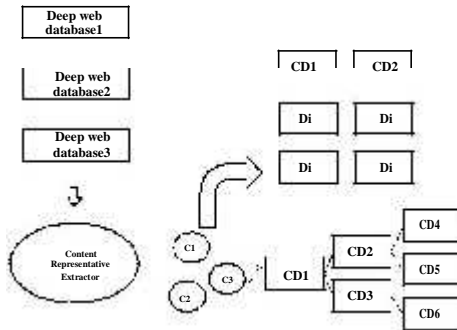


**Figure 1: Deep Web
Classification**

### III.    LIMITATIONS IN EXISTING WORKS

A novel training less approach that automatically classifies deep websites into hierarchical categories of Wikipedia. These categories are merely nodes for organizing the articles in the Wikipedia. The quality of categorized documents is continuously improved by the huge community of Wikipedia authors. The articles in Wikipedia is not strictly categorized i.e. one article may belongs to more than one categories. A huge numbers of categories have been used for identifying documents topic and for measuring the semantic relatedness among text snippets[20], [21].

*A.    Content Representative Extraction*

The limitations in the existing system is that for to classify a deep web database in a domain we need to understand about the kind of content contained in its archives. Thus content representative is that piece of information that tells about the domain of the database. In relation with classification of deep web sources, several techniques have been proposed in the literature to find the domain of online data sources. A strong representative which briefly encounters all the majors of a data source and understandable to a computing resource helps in efficient classification [22] .

Currently two approaches are adopted to extract content representative from the deep web sources through query probing and through visible form features. The first approach peeps inside the deep web database and the second approach tries to find it on the surface of the deep web database source from its attributes.

*B.    Trainingless based classification*

The trainingless based classification involves the process of chosing a relevant keyword and then perform classification but it is not necessary to have for trained set of data within the database. But irrelevant or related words cannot be classified in this case. Further words approach lack semantics, it is not possible to find "general" and "specific concept of a word. Word disambiguation cannot be performed using this approach

### IV.    PROPOSED SYSTEM

In the proposed system, we implement the technique as a web application and the results were evaluated. It was observed that our technique is comparatively better than existing techniques in the following perspectives.

    a)  It is training less which removes the drawback of requiring the representative training data set

    b)  It uses domain independent ontology to semantically classify the deep web database instead of building a category ontology for each and every category domain

    c)  It uses self descriptive Wikipedia category hierarchy keeping in view diverse nature of deep web instead of non-descriptive limited category domains

    d)  It can classify both structure and un-structured deep websites.

In our work we will experiment with semantic deep web database selection methods based on semantic understanding of user''s query based on the prototype of our deep web classification method. Also We Provide a user friendly environment to access and visualize the connection of the input query and the ontology of deep web sources. Unlike the previous system, the proposed model focuses on extracting the pages based on ontology model created among the deep web source based on level and its relationship with each other. We experiment our system with different query sites, which proves our work to be more consistent and reliable when compared to the previous researches, A Novel Algorithm for classifying the ontology and citations, and extraction model is used accommodate the source pages in local storage.

*A.    Web URL Identification*

Identify where the resource is available and the mechanism for extracting it. The uncertainty in usage stems from belonging to different stylistic representation of the semantics of the terms involved.

**TABLE I
Wikipedia Categories for the term "Book" and "List
of Books"**

| Concept term | Wikipedia Classification |
|---|---|
| Single term:Book | Book \| Publishers |
| Comma separated term: List of Books | Environmental books \| Research books \| Journal publications |

*B.    Information Extraction and Parsing*

They involve in accessing HTML as tokens which decodes entities in attributes. Regular expression of the system

success in many web-related areas, such as

- Page ranking in the search engine.
- Google Web page community construction
- Web search improvement
- Web clustering and visualization
- Relevant page finding

processes the spaces and new lines, single and double quotes, HTML comments, and a lot more. The next step up from a regular expression is an HTML tokenizer, HTML Parser to extract information from HTML files. Using these techniques, you can extract information from any HTML file. And automatic passage extraction methods from the body may be worthwhile. Implications of the findings for aids to summarization, and specifically the Text

### C.  Knowledge Source Representation and filtering

Based on the previous extraction process, the knowledge source is constructed and Several external knowledge sources are learned continuously. Many websites are specified for this category. In our approach we employ a domain independent ontology to extract semantics i.e. concepts for the meta-information of the deep web site and the structure and categories to perform the classification task. We represent the meta-information by knowledge source . We use Singular value Decomposition matrix to identify the concepts appearing within the meta-information and apply filters that restrict specific domains and keyword based pages. ontology extracted provides the data in the form are proceeded for further levels. The major purpose of ontology is to relate web content and data to a standard set of subject concepts and provide a fixed set of reference points in a global knowledge space. These subject concepts have defined relationships among them, and can be employed as binding or attachment points for any web content or data and to filter domain specific and content specific informations.

### D.  Ontology Based semantic Classification

Extended Co-Citation also finds some large trustworthy pages present in the website, which provides more accuracy. Please notice that Extended Co-Citation uses no training data, and the testing data is manually analyzed by classifying the depth of the website. Therefore, we believe Extended Co-Citation performs iterative computation to find out the set of ontology and classify the pages accordingly. In order to test its accuracy, we randomly select websites and manually find out their groups. We find the connections and relationship between use the pages and the classification types as the standard fact

### V.    PROBLEM FORMULATION

There are many ways to find relevant pages. Netscape uses web page content analysis, usage pattern information and linkage analysis to find relevant pages.[23] Among the approaches of finding relevant pages, hyperlink analysis has its own advantages. Primarily, the hyperlink is one of the most obvious features of the web and can be easily extracted by parsing the web page codes. Hyperlink analysis has proven

HITS (Hyperlink-Induced Topic Search) algorithm is the existing hyperlink based algorithm which focuses on finding relevant child pages from the page source. Page source is a set of related pages. If the page source is not constructed properly, topic drift problem would arise.

### Definition 1

„Xa‟ denotes the set of keywords that are separated by commas, „Ya‟ denotes the set of single words which are from the meta tag of the html code in the deep web, „Za‟ denotes retrieval of set of subjects from „Ya‟. The Wikipedia categories that are obtained from comma separated word is provided as „Cx‟, „Cy‟ are the set of categories that are obtained from set of words and „Cz‟ are the set of categories that are formed from subject concept.

### Definition 2

The semantic relatedness for the Rya and Rza are found to be Rsr=Rza<string match> Rya ( Ja , Sa , Tr Lr ) AND Rza (Ja , Sa , Tr , Lr ) <string match> Rya ( Ja , Sa , Tr , Lr )Where Ja is the main object, Sa is the subcategories, Tr is the pages and Lr is the parent hierarchy classification of the desired attribute Rya and Rza.

On providing the of input the tokenizer performs the functionalities like removing the stop words and takes only the certain keywords as input. Further mostly the input are found to be comma separated words, so the tokenizer returns all the comma separated words too.

So in the proposed system we have introduced the process of finding the concepts which is found to be in the existing material that are subject-related and theoretical rather than physical. After identifying the concept the related input is mapped with Wikipedia hierarchy. Now the identified list of categories are mapped to the URL search source and the database that performs knowledge management. Thus on mapping the URL provides the knowledge database with deep web sites.

**Alogorithm**
**Input**
   The data model that are extracted from the HTML code of the deep website.

**Output**
   Mapping the input with the hierarchy of the Wikipedia category

1. X,Y:=Input
2. Z:=Concept identification
3. $C_x$, $C_y$, $C_z$ :=Wikicategories
4. $R_{sr}=R_{za}$<string match> $R_{ya}$ ( $J_a$ , $S_a$ , $T_r$ , $L_r$ ) AND $R_{za}$ ($J_a$ , $S_a$ , $T_r$ , $L_r$ ) <string match> $R_{ya}$ ( $J_a$ $S_a$ , $T_r$ , $L_r$ )

## VI. EXPERIMENT RESULTS

According to our Proposed model, data instances of the same type of pages have the same path from the root in the of the input parent page. Thus, our algorithm does not need to

merge similar page ontology from different levels and the task to merge multiple pages can be broken down from a tree level to a string level of same parent page or different.

#### A. Root Node Selection

The process of root node selection involves providing the specified URL that is needed to be analyzed. Starting from root node of all input web page is analyzed for ontological structure based on semantic relationship is constructed, with a filtering approach is introduced to find and remove the inappropriate contents from the entire web page, which is depth invariant, additionally to keyword based filtering, the system applies a domain extension based filtering and file extension based filtering to prevent malicious attacks and other attacks.

**Table 2: URL selection**

| Input type | URL choosen |
|---|---|
| News | www.msn.com |
| Movies | www.movie.com |
| Books | www.books.com |

#### B. Related Web links

The system discovers a methodology to find the non working pages present in any depth is detected and classified which helps in detecting the unwanted and broken pages in the entire web structure, our algorithm applies a new multiple Co-Citation Approach to their initial to each and n-level child nodes. There are multiple advantages in the proposed design.

**Table 3: Related sites**

| Domain | Description | Sites |
|---|---|---|
| Wikipedia | Wikibooks, Wikimedia,Wikiversity | 313 |
| News | National,International, Weekend, Sports, Entertainment, Video | 256 |
| Books | Book stores, Book Authors, Book title | 271 |

The number of child nodes under a parent node is much smaller than the number of nodes in the whole depth or the number of HTML tags in a Webpage, thus, the effort for multiple depth level analysis is made which is much efficient than the other deep page classification methods.

#### C. Filtering

The nodes with the same tag name can be better differentiated by the subtrees they represent, which is an important feature not used in previous experiments.

**Table 4: Category coverage**

| Domain | Root Category | Child links | Dataset coverage |
|---|---|---|---|
| News | Report National | 136 | 97% |
| Movies | Movie type | 98 | 95% |
| Books | Title, Author | 154 | 99% |

Instead, our algorithm will recognize such nodes as peer nodes and denote the same symbol for those child nodes to facilitate the experimental process successful for news websites(ndtv, bbc, cnnibn ), commercial website (yahoo, msn, aol, indiatimes), Video websites( youtube), business websites (jabong, ebay, snapdeal, flipcart) directory site (justdail, suleka, a2znational) also for wiki website.The filtering process provides in removing hanging tags and generate a newer classification.
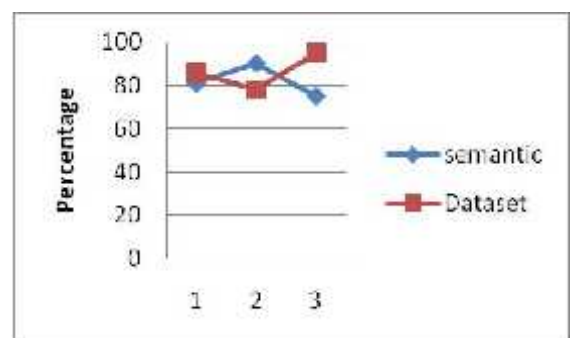
#### D. Semantic Ontology

This table provides the semantic ontology for the domains.

**Table 5: Dataset damains**

| Domain | Semantic | Dataset coverage |
|---|---|---|
| News | 81 | 86% |
| Movies | 90 | 78% |
| Books | 75 | 95% |

**Figure 2: Comparative analysis**

current system and a hyperlink analysis is made to track the continuity of data piracy and instead of having page level web data extraction, the system can be implemented for web site level information analysis for template pages and comparison.

\

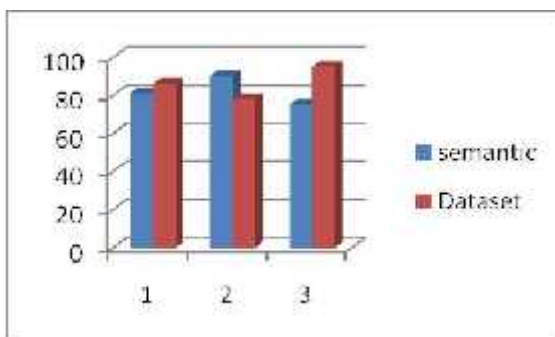The measure of semantic and dataset are given



**Figure 3: Measuring semantic deep web classification**

Thus the results are analyzed for the classification of the deep web by semantic method.

### VII.   CONCLUSION

According to our Proposed model, data instances of the same type of pages have the same path from the root in the of the input parent page. Thus, our algorithm does not need to merge similar page ontology from different levels and the task to merge multiple pages can be broken down from a tree level to a string level of same parent page or different . Starting from root nodes of all input web page is analyzed for ontological structure based on semantic relationship is constructed, with a filtering approach is introduced to find and remove the inappropriate contents from the entire web page, which is depth invariant, additionally to keyword based filtering, the system applies a domain extension based filtering and file extension based filtering to prevent malicious attacks and other attacks. The system discovers a methodology to find the non working pages present in any depth is detected and classified which helps in detecting the unwanted and broken pages in the entire web structure, our algorithm applies a new multiple Co-Citation Approach to their initial to each and n-level child nodes. The future systems can be enhanced to extract web pages and sites from web server which is not implemented with the

REFERENCES

[1] The Deep Web: Surfacing hidden value. Accessible at http://brightplanet.com, 2000.

[2] Madhavan, J., Cohen, S., Dong, X. L., Halevy, A. Y., Jeffery S. R., Ko, D. and Yu, C. Web scale data integration: You can afford to pay as you go. In the proceedings of Conference on Innovative Data Systems Research (CIDR),2007, 342–350.

[3] Chang, K. C. C., He, B., Li, C., Patel, M. and Zhang,Z. Structured databases on the web: Observations and Implications. In the proceedings of International Conference on Management of Data (ACM SIGMOD), 2004, 61–70.

[4] Madhavan, J., Afanasiev, L., Antova L. and Halevy, A.Y. Harnessing the Deep Web: Present and Future. In the proceedings of Conference on Innovative Data Systems Research (CIDR), 2009.

[5] Barbosa, L. and Freire, J. Searching for Hidden-Web Databases. In the proceedings of International Workshop on the. Web and Databases (WebDB), 2005, 1–6.

[6] Madhavan, J., Ko, D., Kot, L., Ganapathy, V., Rasmussen A. and Halevy, A. Y. Google's Deep Web crawl. Incthe proceedings of International Conference on Very Large Database Endowment (PVLDB), 2008, 1241-1252.

[7] Raghavan, S. and Garcia-Molina, H. Crawling the Hidden Web. In the proceedings of International conference on Very Large Databases (VLDB), 2001, 129–138.

[8] Xian, X., Zhao, P., Fang, W., Xin, J. and Cui, Z. Automatic Classification of Deep Web Databases with Simple Query Interfaces. In the proceedings of International Conference on Industrial Machatronics and Automation (ICIMA), 2009, 85–88.

[9] Ipeirotis, P. G., Gravano, L. and Sahami, M. Automatic Classification of Text Databases through Query Probing. In ACM SIGMOD workshop on the Web and Databases (WebDB), 2000, 245–255.

[10] Xu, H., Hau, X., Wang, S. and Hu, Y. A method of Deep Web Classification. In the proceedings of International Conference on Machine Learning and Cybernetics (ICMLC), 2007, 4009–4014.

[11] Nie, T., Shen, D., Yu, G. and Kou, Y. Subject-

Oriented Classification Based on Scale Probing in the Deep Web. In the proceedings of International Conference on Web-age Information Management (WAIM), 2008, 224– 229.

[12] Lin, P., Du, Y., Tan, X. and Lv, C. Research on Automatic Classification for Deep Web Query Interfaces. In International Symposium on Information Processing (ISIP), 2008, 313–317.

[13] Le, H. and Conrad, S. Classifying Structured Web Sources Using Support Vector Machine and Aggressive Feature Selection. In Lecture Notes in Business Information Processing, Vol. 45, 2010, 270–282.

[14] Zhao, P., Huang, L., Fang, W. and Cui, Z. Organizing Structured Deep Web by Clustering Query Interfaces Link Graph. In Lecture Notes in Computer Science, Vol. 5139, 2008, 683–690.

[15] Barbosa, L., Freire, J. and Silva, A. Organizing hidden-web databases by clustering visible web documents. In the proceedings of International Conference on Data Engineering (ICDE), 2007, 326–335.

[16] He, B. Tao, T. and K. Chang, K. C. C. Organizing structured web sources by query schemas: a clustering approach. In the proceedings of Conference on Information and Knowledge Management (CIKM), 2004, 22–31

[17] Su, W., Wang, J. and Lochovsky, F. Automatic Hierarchical Classification of Structured Deep Web Databases. In the proceedings of International Conference on Web Information System Engineering (WISE), 2006, 210– 221.

[18] Medelyan, O., Milne, D., Legg, C. and Witten, I. Mining meaning from Wikipedia. In International Journal of Human-Computer Interactions (IJHCI), Vol. 67(9), 2009, 716–754.

[19] Noor, U., Rashid, Z. and Rauf, A. A survey of automatic deep web classification techniques. In International Journal of Computer Applications (IJCA), Vol. 19(6), 2011, 43–50.

[20]Schonhofen, P. Identifying document topics using the Wikipedia category network. In the proceedings of International Conference on Web Intelligence (IEEE/WIC/ACM), 2006, 456–462.

[21] Huynh, D., Cao, T., Pham, P. and Hoang, T.Using Hyperlink Texts to Improve Quality of IdentifyingDocument Topics Based on Wikipedia. In the proceedings of International Conference on Knowledge and Systems Engineering (ICKSE), 2009, 249–254.

[22] Noor, U., Rashid, Z. and Rauf, TODWEB: Training-lessOntology based Deep Web Source Classification ,2011 ACM978-1-4503-0784 iiWAS''11, 2011,5–7.[23] Jingyu Hou, Yanchun Zhang,Effectively Finding Relevant Web Pages from Linkage Information, IEEE Computer Society, 2003 (vol. 15 no. 4),pp. 940.

The Author, Sreekrishna.M, from Coimbatore born on 12th December 1989 is a final year student doing Master of Engineering in Computer Science at Angel College of Engineering and Technology,Tirupur and has a Bachelor of Engineering degree in Computer Science and Technology from Ponjesly College of Engineering, Nagercoil, earned the degree in the year 2011. Her Research area is Data Mining and have attended PSG ACM Workshop for her future research work. She presented papers in national and International conferences.

The Author, Sundaramoorthy.P, pursuing his PhD degree at Anna University, Chennai in the domain of Cloud Computing. Now he is working as a Assistant Professor in the Department of CSE at Angel College of Engineering and Technology, Tirupur. He has presented his papers in National and International conferences in various domains and have attended workshops related to his work. His area of interest is Cloud Computing and Data Mining.

The Author, Dr.T.Rajendran completed his PhD degree in 2012 at Anna University, Chennai in the Department of Information and Communication Engineering Tamilnadu, India. Now he is working as a Dean for Department of CSE & IT at Angel College of Engineering and Technology, Tirupur. His research interest includes Distributed Systems, Web Services, Network Security and Web Technology. He is a life member of ISTE & CSI. He has published more than 51 articles in International/ National Journals/Conferences. He has visited Dhurakij Pundit University in Thailand for presenting his research paper in International conference. He was honored with Best Professor Award 2012 by ASDF Global Awards 2012, Pondicherry.