

Multi-Stage Web Crawler for Deep Web Resources

Patil Ashwini M.¹, Lambhate Poonam D²

^{1,2} Department of Computer Engineering

^{1,2} JayantraoSawant College of Engineering, Hadapsar Pune-28, Savitribai Phule Pune University, Pune, India

Abstract- Surface web is often websites of companies, peoples and bloggers. The opposite term to the deep web is the surface web. The deep web contents are not indexed by the standard search engine. Also the contents of deep web changes rapidly. Thus locating the contents of deep web effectively is difficult. To address the problem we propose a web crawler having two stages. In the first stage, crawler will retrieve the relevant sites for a focused crawler and assign priorities to the web sites according to their relevance. In the second stage, the crawler does in-site searching using adaptive link learning which finds searchable forms. The design of link tree data structure helps to avoid missing the relevant searchable forms.

Keywords- Meta Search, Web crawler, Page Ranking, Reverse searching

I. INTRODUCTION

The deep web is also called as invisible web or hidden web. The contents of deep web are not indexed by standard search engines. The opposite term to the deep web is the surface web.

Surface web is the thought of static websites (though connected to deep web databases, such as Amazon.com) Examples of surface web pages are Google, Facebook.

Surface web is often websites of companies, peoples and bloggers.

Deep web have images of court records, maybe archives of old newspapers. The deep web is largely academic databases and government archives which can not be seen by surface web. The image below shows the meaning of surface web versus deep web.

The iceberg shows the deep web and the surface web.



Fig: Iceberg showing deep web and surface web.

A crawler visits Web sites and retrieves their pages and other information for later processing by a search. The major search engines on the Web all have a crawler, also called as a "spider" or a "bot".

Crawlers crawl through a website one page at a time until all pages on the site have been read, it will follow the links to other pages and thus called a crawler.

Web crawlers are mainly used to create a copy of all the visited pages, engine in order to create entries for a search engine index for effective search.

Crawlers can be used for checking links or validating HTML code. Also, crawlers can be used to collect specific types of information from Web pages, such as obtaining a large number of e-mail addresses (usually for spam).

Berkeley estimated that the deep web contains approximately 91,850 terabytes and the surface web is only about 167 terabytes in 2003[7].

The deep web is about 500-550 times larger than the surface web[2]. The deep web may contain valuable information, and keep constantly changing[1],[2]. To locate deep web databases is a challenging task as they are not registered with any search engine. Thus there arises a need for a web crawler which can explore deep web databases accurately and quickly.

The rest of the paper is organized as follows. Section II is about Literature Review. Third section has System Architecture. The fourth section describes the System Analysis. Fifth section explains the algorithm. Experimental results are presented in section six.. Concluding remarks are given in section seven.

II. LITERATURE REVIEW

The different techniques were proposed to locate deep web efficiently.

In [14], author survey the part of the deep Web consisting of dynamic pages in one particular national domain. The approximate calculation of the national deep Web is performed using the proposed sampling techniques. The data

from the Russian search engine “yandex” and Russian segment called “Runet” is used. Two sampling techniques are used to estimate the number of deep web sites.

- i) Random Sampling Technique
- ii) Stratified Random Sampling Technique.

It is seen that the proportion of deep web is high in highly sited web sites than less cited web sites. The problem with Generic Crawlers is that they fetch all searchable forms and can not focus on specific topic.

In [13], there are two ways to access the deep web. The first is to create vertical search engines for specific domains and the second way is surfacing. The prototype system for surfacing Deep-Web content is proposed. The proposed algorithm efficiently traverses the search space and identifies the URLs that can be indexed by the Google search engine.

The deep web may contain data of two types. The first is text based and the second is the structured entities. The example for entity oriented web is online shopping sites. The techniques used before are not effective to access entity oriented deep web. Thus in [6], the prototype system to crawl entity oriented deep web is proposed with the techniques including query generation, empty page filtering and the URL deduplication. But in template generation, the parsing handles only “GET” forms but not the “POST” forms. The proposed techniques are seen useful for crawling entity oriented deep web.

In [18], apprentice is used to help crawler to assign priorities to URLs in crawl frontier using predefined set of features and events related to crawler. The apprentice takes online lessons from the focused crawler. When the apprentice get ready with enough examples, the crawler takes suggestions from apprentice while crawling to better prioritize unvisited URLs in crawl frontier. This achieves higher rate of retrieval of relevant pages. The online relevance feedback helps to reduce false positives.

In [8], surface web is the HTML pages interconnected with hyperlinks. Hidden web that is the dataset of an organization, while accessing through interface, only few tuples are retrieved and thus the search engines can not effectively crawl the hidden web. The authors solved the problem for crawling hidden web when data set is numeric, categorical or both through proposed algorithms.

In [16], the proposed system focus on a specific topic reducing the need to crawl a large number of unrelated pages

but covers a broad area for search. The stopping criteria helps to avoid visiting unproductive links within a site. The framework proposed can also be used to build focused crawler for different domains.

In [15], developed a new framework called ACHE (Adaptive Crawler for Hidden Web Entries). The framework removes the limitations of form focused crawler. The crawler learns patterns adaptively. The adaptive crawler gives high quality results in less time.

III. SYSTEM ARCHITECTURE

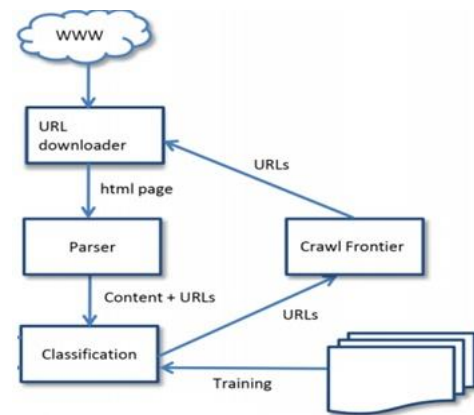


Fig: System Architecture

IV. SYSTEM ANALYSIS

In the system, the URLs are first fetched from search engine. The clustering is done on retrieve pages as relevant or irrelevant. After clustering according to priority the pages are displayed to the user. Then the user may go for the deep search as per the requirement.

V. SYSTEM ALGORITHM

Step1. Retrieve Urls From Google Search Engine.

Step2. Compare The Keyword In The Query With The Description And Title Of The Url.

Step3. Clustering Is Done Of Fetched Urls.

Step4. Ranking Is Assigned To Pages Using Cosine Similarity.

Step5. The Relevant Urls Will Be Displayed To The User According To Their Priority.

Step6. The User Can Go For Deep Search Through The Links Displayed.

VI. EXPERIMENT AND RESULT

The expected results using the web crawler are shown using graphs. It shows that Our web crawler will retrieve the maximum number of searchable forms for different domains.

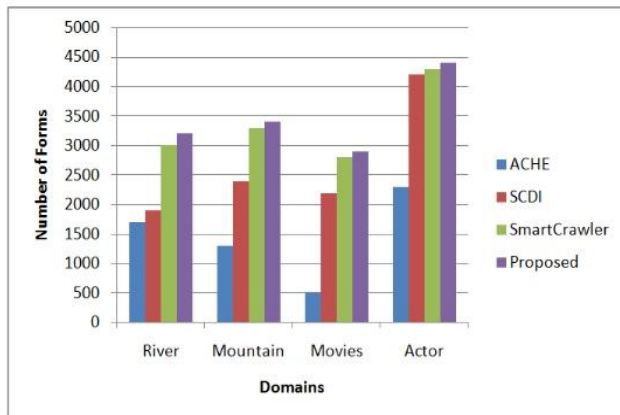


Fig: Expected Results

VII. CONCLUSION

Literature survey and analysis of some of the important deep web crawlers was carried out to find their advantages and limitations. A comparative analysis of existing deep web crawlers was also carried out on the basis of various parameters and it is concluded that a new architecture for deep web crawler was required for efficient searching of the deep web information by minimizing the limitations of the existing deep web crawlers. Hence to improve searching efficiency, a new crawler architecture is proposed having two stages, one for retrieving the relevant sites and the second will help for deep search.

ACKNOWLEDGMENT

I would like to thank my guide Prof .Poonam D.Lambhate for her help and guidance throughout this project and the semester, without her this would not have been possible.

REFERENCES

- [1] Poonam P. doshi, dr. emmanuel m : feature extraction techniques using semantic based crawler for search engine. In. proceedings of international conference on computing, communication and energy systems. 2016.
- [2] Feng Zheo, Jingyu Zhou, Chang Nie, Heqing Huang, Hai Jin. SmartCrawler: A Two-Stage Crawler for Efficiently Harvesting Deep-Web Interfaces. IEEE Transactions on Services Computing, vol 99, 2015
- [3] Booksinprint. Books in print and global books in print access. <http://booksinprint.com/>, 2015.
- [4] Idc worldwide predictions 2014: Battles for dominance – and survival – on the 3rd platform. <http://www.idc.com/research/Predictions14/index.jsp>, 2014.
- [5] Infomine. UC Riverside library. <http://lib-www.ucr.edu/>, 2014.
- [6] Yeye He, Dong Xin, Venkatesh Ganti, Shriram Rajaraman, and Nirav Shah “crawling Deep Web Entity Pages” In Proceedings of the sixth international conference on web search and data mining, pages 355-364, ACM, 2013
- [7] Martin Hilbert. How much information is there in the “information society”? Significance, 9(4):8–12, 2012.
- [8] Cheng Sheng, Nan Zhang, Yufei Tao and Xin Jin”Optimal Algorithms for Crawling a Hidden Database in the Web”. Proceedings of the VLDB Endowment,5(11):1112-1123, 2012
- [9] Balakrishnan Raju and KambhampatiSubbarao. Sourcerank: Relevance and trust assessment for deep web sources based on inter-source agreement. In Proceedings of the 20th international conference on World Wide Web, pages 227–236, 2011.
- [10] Denis Shestakov. Databases on the web: national web domain survey. In Proceedings of the 15th Symposium on International Database Engineering & Applications, pages 179–184. ACM, 2011.
- [11] Olston Christopher and Najork Marc. Web crawling. Foundations and Trends in Information Retrieval, 4(3):175–246, 2010.
- [12] Denis Shestakov and Tapio Salakoski. Host-ip clustering technique for deep web characterization. In Proceedings of the 12th International Asia-Pacific Web Conference IEEE 2010
- [13] Jayant Madhavan, David Ko, Lucja Kot, Vignesh Ganapathy, Alex Rasmuseen and Alon Halevy “Google’s Deep Web Crawl” Proceedings of the VLDB Endowment, 1(2):1241-1252, 2008
- [14] Denis Shetakov and Tapio Salskoski “on estimating the scale of national deep web” In Database and Expert systems Applications, pages 780-789. Springer, 2007

- [15] Luciano Barbosa and Juliana Freire "An Adaptive Crawler for Locating Hidden-Web Entry Points" In Proceedings of the 16th international conference on world wide web, pages 441-450. ACM, 2007
- [16] Luciano Barbosa and Juliana Freire "Searching for Hidden-Web Databases" In WebDB, Pages 1-6, 2005
- [17] Peter Lyman and Hal R. Varian. "How much information?" 2003. Technical Report, UC Berkeley, 2003.
- [18] Soumen Chakrabarti, Kunal Punera and Mallela Subramanyam "Accelerated Focused Crawling through Online Relevance Feedback" In Proceedings of the 11th international conference on world wide web, pages 148-159, 2002.