



International Journal of Innovative Research in Computer and Communication Engineering

(A High Impact Factor, Monthly, Peer Reviewed Journal)

Vol. 4, Issue 1, January 2016

Smart Crawler: A Two Stage Crawler for Efficiency Harvesting Deep Web Interfaces

Satyawan Dongare, Komal Gawali, Minal Pathak, Prof Bharati Gaikwad.

Department of Computer Engineering, G.S. Moze College of Engineering, Balewadi, Pune, MH, India

ABSTRACT: Due to heavy usage of internet large amount of diverse data is spread over it which provides access to particular data or to search most relevant data. It is very challenging for search engine to fetch relevant data as per user's need and which consumes more time. So, to reduce large amount of time spend on searching most relevant data we proposed the "Advanced crawler". In this proposed approach, results collected from different web search engines to achieve meta search approach. Multiple search engine for the user query and aggregate those result in one single space and then performing two stages crawling on that data or Urls. In which the sight locating and in-site exploring is done f or achieving most relevant site with the help of page ranking and reverse searching techniques. This system also works online and offline manner.

KEYWORDS: Meta search, Two stage crawler, Page Ranking, Reverse searching.

I. INTRODUCTION

Internet is important part of our day to day life. It is an indivisible part of modern generation as well as old generation. To get answer of most common question there is a need of an Internet, and the Internet gives the birth to the WWW and there are huge amount of data spread over WWW. There are many Search Engine are used over the WWW, but the mostly used search engine are Google, yahoo, msn. In the race of searching they keep their stamp because their precision rate and internal algorithm but the problem with these general search engine is that they are best for surface web searching but not too good for deep web searching in which what an end user expecting that maximum relevant documents must retrieved with his query. But as the space on WWW is increasing it contains a vast amount of data, of an valuable information and these information cannot access properly by web indices in web search engine (e.g. Google, Baidu) then to overcome these problem there is need of efficient harvesting which accurately and quickly explore the deep web, it is challenging to locate a deep web database because they are not register with any search engine. To address this problem previously working on two types of crawler first is generic crawler and second is focused crawler, Focused crawler search automatically on-line database from to search engine, generic crawler is hidden or adaptive crawler [1].

So to harvest the deep web and to provide the answer of user query with minimum effort we are proposing the Advance crawling concept which is based on MetaSearch strategies and smart two stage crawling [1]. These Search engine gives the answer of basic question but the need of corporate world is increased day by day and they need answer of harder question which is unanswerable. [9].The World Wide Web is huge repository of the information, complexity is more when to accessing a data from to WWW i.e. there is need of efficient searching technique to extract appropriate information from the web. A Meta search engine is a search engine tool that send user request to several search engines concurrently and the aggregates the result into single list and displays them according to the relevance. In this approach meta search engine enables user to enter search query once and access several search engine simultaneously in this strategy advantage is that maximum relevant documents can be retrieved but condition is that those retrieved documents must satisfy the threshold value which is a boundary conditions [2]; In this approach critical task is to combine several search engine with proper ranking of relevant documents.



International Journal of Innovative Research in Computer and Communication Engineering

(A High Impact Factor, Monthly, Peer Reviewed Journal)

Vol. 4, Issue 1, January 2016

II. LITERATURE SURVEY

Around 1993, ALIWEB was grown up as the web page identical to Archie and Veronica. Instead of cataloguing files or text records, webmasters would submit a special systematized file with site information[4]. The next development in cataloguing the web came later in 1993 with spiders. Like robots, spiders scoured the web for web page information. These early versions looked at the titles of the web pages, the header information, and the URL as a source for key words. The database techniques used by these early search engines were primitive. For example, a search process would give up hits (List of Links) in the order that the hits (List of Links) were in the database. Only one of these search engines made effort to rank the hits (List of Links) according to the website's relationships to the key words. The first popular search engine, Excite, has its roots in these early days of web Classifying. The Excite project was begun by a group of Stanford undergraduates. It was released for general use in 1994[4]. From the recent few years there are many methodologies and techniques are proposed searching the deep web or searching the hidden data from one of the www. The design of Meta search engine and same deep web searching site are the example of this engine. The first Meta search engine was created during the year of 1991-1994. It provides the access to many search engines at a time by providing single query as input GUI. And its name as a Meta crawler Proposed in university of Washington [3]. And after onward the work is still going on to save the data from diving in Deep Ocean of internet there are one of the best example of this invention is guided Google proposed by cheek Hong dm and rajkumar bagga in which the use Google API for searching and controlling search of Google. The inbuilt method and function library is guided [Google]. In year 2011 web service architecture for Meta search engine was proposed by K SHRINIVAS, P V S SHRINIVAS, A GOVARDHAN according to their study the Meta search engine can be classified into two types first is general purpose search engine and second special purpose Meta search engine. The previous search engine are focused on searching the complete web but year after year to reduce the complexity focus is a to search information in particular domain [2]. Information retrieval is a technique of searching and retrieving the relevant information from the database. The efficiency of searching is measure using precision and recall .Precision Specifies the document which are retrieve that are relevant and Recall Specifies that whether all the document that are retrieve are relevant or not. Web Searching is also type of information retrieval because the user searchers the information on web. The information that is search on web is called as web mining. Web mining can be classified in three different types Web Content Mining, Web Structure Mining, Web Usage Mining. To retrieve the complex query information is still checking for search engine is known as deep web. Deep web is invisible web consist of publicly accessible pages with information in database such as Catalogues and reference that a not index by search engine [2]. The Deep web is rapidly growth day over day and to locate them efficiently there is need of effective techniques to achieve best result. Such a system is effectively implemented is Smart Crawler, which is a Two Stage Crawler efficiently harvesting deep web interfaces. By using some basic concept of search engine strategies they achieve the good result in searching of most significant data. Those data techniques are as reverse searching, incremental searching.

III. RELATED WORK / BACKGROUND

A. Internet archive Crawler

In 1997, Mike Burner designed the Internet Archive Crawler[2] was the first paper that focused on the challenges caused by the scale of web. It uses multiple machine to crawl the web and it crawl on 100 million URLs[1]. Each crawler process read a list of seed URLs for its assigned sites from disk into per-site queue, and then it uses asynchronous I/O instructions to fetch pages from these queues in parallel. It has also deal with the problem of changing DNS records, so it keeps the historical archive of hostname to IP mapping.

B. Google Crawler

Later in 1998, The original Google crawling system consist of a five crawling components which was running in various process and download the pages[2]. Each crawler process used asynchronous I/O instructions to fetch the data from up to 300 web servers in parallel. Then all the crawlers transmit downloaded pages to a single Store Server process that compressed the page and store them on disk[1]. Google Crawler was based on C++ and Python. This



International Journal of Innovative Research in Computer and Communication Engineering

(A High Impact Factor, Monthly, Peer Reviewed Journal)

Vol. 4, Issue 1, January 2016

crawler was integrated with the indexing process(text parsing was done for full-text indexing and also for URL extraction).

C. Mercator Web Crawler

Heydon and Najork present a web crawler which was highly scalable and easily extensible [3][1]. It was written in Java. The first version was non-distributed and later the distributed version was made available which split up the URL space over the crawlers according to host name and avoid the potential bottleneck of a centralized URL server.

D. Web Fountain crawler

In 2001, another distributed and modular crawler represented by IBM[4][1]. It has three major component, Multi threaded crawling processes, duplicate content and central controlled process responsible for assigning work. It was written in C++ and used MPI to facilitate the communication between the various process. It was deployed on a cluster of 48 crawling machine.

IV. PROPOSED SYSTEM

In this proposed system we use a meta search engine; In meta search engine the searching result is accumulated by using multiple search engines. Actually it is good in finding the unique key word phrases, quotes, and Knowledge encompasses in the full text of web pages. And Search engines allow user to enter keywords and then examine this keyword in its table followed by database. We propose a novel two-stage framework to address the problem of searching for hidden-web Resources. Our site locating technique employs a reverse searching strategy(e.g. By using Google's link: facility to get pages guiding to a given link) and incremental two-level site prioritizing technique to discover more relevant sites, achieving more data sources. During the in-site exploring stage, we have constructed a link tree for balanced link prioritizing, eliminating bias toward web pages in popular directories. We introduce an adaptive learning algorithm that performs online feature selection and uses these features to automatically create link rankers. In site locating stage, high relevant sites are arranged and the crawling is focused on a topic using the contents of the seed page of sites, achieving more accurate results. During the insite exploring stage, relevant links are arranged for fast in-site searching.

To efficiently and effectively discover deep web data sources, Web Crawler is designed with two stage architecture, site locating and in-site exploring, as shown in architecture diagram. The first site locating stage finds the most relevant site for a given topic, and then the second in-site exploring stage uncovers searchable forms from the site. Specifically, the site locating stage starts with a seed set of sites in a site database. Seeds sites are candidate sites given for web crawler to start crawling, which begins by following URLs from chosen seed sites to explore other pages and other domains. When the number of unvisited URLs in the database is less than a threshold during the crawling process, Then "reverse searching " is perform by web crawler of known deep web sites for center pages (highly ranked pages that have many links to other domains) and feeds these pages back to the site database. Site Frontier fetches homepage URLs from the site databases, which are ranked by Site Ranker to prioritize highly relevant sites. The Site Ranker is improved during crawling by an Adaptive Site Learner, which adaptively learns from features of deep-web sites (web sites containing one or more searchable forms) found. To achieve more accurate results for a focused crawl, Site Classifier categorizes URLs into relevant or irrelevant for a given topic according to the homepage content.

V. CONCLUSION

In this paper we have identified the different kind of general searching technique and Meta search engine strategy and by using this we have proposed an effective way of searching most relevant data from hidden web. In this we are combining Multiple search engine and two stage crawler for harvesting most relevant site. By using page ranking on collected sites and by focusing on a topic, advanced crawler achieves more accurate results. The two stage crawling performing site locating and in-site exploration on the site collected by Meta crawler.



International Journal of Innovative Research in Computer and Communication Engineering

(A High Impact Factor, Monthly, Peer Reviewed Journal)

Vol. 4, Issue 1, January 2016

VI. FUTURE WORK

For the future enhancement we can implement the given system on cross search engine platform like google, yahoo etc. When system will work on cross domain platform we can retrieve whole data as per the user requirements with minimum cost.

REFERENCES

- [1] Olston and M. Najork, "Web Crawling", Foundations and Trends in Information Retrieval, vol. 4, No. 3, pp. 175–246, 2010
- [2] M. Burner, "Crawling towards Eternity: Building an Archive of the World Wide Web," Web Techniques Magazine, vol. 2, pp. 37-40, 1997.
- [3] Allan Heydon and Marc Najork. Mercator: A scalable, extensible web crawler. World Wide Web Conference, 2(4):219–229, April 1999.
- [4] Jenny Edwards, Kevin S. McCurley, and John A. Tomlin. An adaptive model for optimizing performance of an incremental web crawler. In Proceedings of the Tenth Conference on World Wide Web, pages 106–113, Hong Kong, May 2001. Elsevier Science.
- [5] Luciano Barbosa and Juliana Freire. Searching for hidden-web databases. In WebDB, pages 1–6, 2005.
- [6] 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.
- [7] Soumen Chakrabarti, Martin Van den Berg, and Byron Dom. Focused crawling: a new approach to topic-specific web resource discovery. Computer Networks, 31(11):1623–1640, 1999.
- [8] Kevin Chen-Chuan Chang, Bin He, and Zhen Zhang. Toward large scale integration: Building a metaquerier over databases on the web. In CIDR, pages 44–55, 2005.
- [9] 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.
- [10] Denis Shestakov and Tapio Salakoski. +Host-ip clustering technique for deep web characterization. In *Proceedings of the 12th International Asia-Pacific Web Conference (APWEB)*, pages 378–380. IEEE, 2010.
- [11] Denis Shestakov and Tapio Salakoski. On estimating the scale of national deep web. In *Database and Expert Systems Applications*, pages 780–789. Springer, 2007.
- [12] Michael K. Bergman. White paper: The deep web: Surfacing hidden value. *Journal of electronic publishing*, 7(1), 2001.
- [13] Shestakov Denis. On building a search interface discovery system. In *Proceedings of the 2nd international conference on Resource discovery*, pages 81–93, Lyon France, 2010. Springer.
- [14] Bright planet's searchable database directory. <http://www.completeplanet.com/>, 2013.
- [15] Y. Wang, T. Peng, W. Zhu, "Schema extraction of Deep Web Query Interface", *IEEE Transaction On Web Information Systems and Mining*, WISM International Conference 2009.