

Enhance Crawler For Efficiently Harvesting Deep Web Interfaces

Sujata R. Gutte

M.E. CSE Dept

M. S. Bidwe Engineering College,

Latur , India

e-mail: omgutte22@gmail.com

Shubhangi S. Gujar

M.E. CSE Dept

M. S. Bidwe Engineering College,

Latur , India

e-mail: shubhangigujars@gmail.com

Abstract— Scenario in web is varying quickly and size of web resources is rising, efficiency has become a challenging problem for crawling such data. The hidden web content is the data that cannot be indexed by search engines as they always stay behind searchable web interfaces. The proposed system purposes to develop a framework for focused crawler for efficient gathering hidden web interfaces. Firstly Crawler performs site-based searching for getting center pages with the help of web search tools to avoid from visiting additional number of pages. To get more specific results for a focused crawler, projected crawler ranks websites by giving high priority to more related ones for a given search. Crawler accomplishes fast in-site searching via watching for more relevant links with an adaptive link ranking. Here we have incorporated spell checker for giving correct input and apply reverse searching with incremental site prioritizing for wide-ranging coverage of hidden web sites.

Keywords- *Focused Crawler, Hidden Web , Feature Selection, Ranking;*

I. INTRODUCTION

Net search tools attempt to get data as relevant as possible. One of the parts of net search tools is the Web Crawler. A web crawler is a program that goes across the net assemble and collect data in a database for further analysis. The method of web crawling includes gathering pages from the net and organizing them out in a manner that the search engine can use them efficiently.

The objective is to do it fast and efficiently without much barrier with the operation of the remote server. A web crawler starts with a URL also called as seed. The crawler visits the links in the list and it also looks for hyperlinks to additional web pages. It then adds them to the present list of URLs in the list. This procedure of crawler visiting URLs depend on upon set of rules for the crawler. Usually, crawler crawls the URLs in the present list incrementally. In addition to that, crawler gather data from the Web Pages. This paper is organized into following sections. Section II describes the investigates done on focused crawlers earlier. Section III states the goal of proposed architecture for focused crawler. in section IV Methodology, algorithms are described. Initial result obtained is show in section V.

II. LITERATURE SURVEY

Form-Focused Crawler (FFC) [6] and Adaptive Crawler for Hidden-web Entries (ACHE) [7] mechanically search online databases for a defined topic. FFC includes link, page and form classifiers for focused crawling of web forms and ACHE is an improvement of FFC with modules for form filtering and adaptive link learner. The FFC and ACHE are focused crawlers planned for searching concerned deep web interfaces. FFC Crawler performs a broad search by focusing the search on definite topic efficiently by learning to detect promising links. The crawler uses two type of classifiers the page and the link classifier to conduct the search. And the third classifier, the form classifier filters out unusable forms. The page The hidden

web content is the data that cannot be indexed by search engines as they usually stay behind searchable web interfaces [2]. A Crawler encounters a variety of web pages during a crawling procedure . For efficient crawling and wide-ranging coverage, ranking and prioritizing links of dissimilar sites is necessary. In earlier work two types of crawlers were proposed, these crawlers are generic crawlers and focused crawlers. Generic crawlers are mostly created for representing hidden web and construction of directory for hidden web resources. The search is not limited to a specific topic, but tries to get all searchable forms [3,4]. Thus, generic crawlers gathers all searchable forms and does not focus on a exact topic. Focused Crawler is a web crawler for fetching web pages that are related to a particular area of interest [5]. It gathers the documents that are related to a given subject. It is called as a Topic Crawler as the outcome of the way it works. The focused crawler decides the significance of the document before crawling the page. It predicts if the given page is relevant to a particular topic and how to proceeds. The essential point of preference of this sort of crawler is that it needs less tools assets.

Form-Focused Crawler (FFC) [6] and Adaptive Crawler for Hidden-web Entries (ACHE) [7] automatically search online databases for input topic. FFC includes link, page and form classifiers for focused crawling of web forms and ACHE is an improvement of FFC with works for form filtering and adaptive link learner. The FFC and ACHE are focused crawlers intended for searching interested hidden web interfaces. FFC Crawler performs a broad search by focusing the search on specified topic well by learning to identify promising links. The crawler uses two classifiers the page and the link classifier to perform the search. And the third classifier, which is the form classifier filters out useless forms. The page classifier is trained so as to categorize pages according to topics specified. When the crawler fetches a page links are extracted from homepage. The link classifier examines links which are extracted from topic specific pages and adds them to the crawler frontier according to their precedence. The link classifiers trained to identify links that can to lead to pages which consist of searchable form interfaces. But this approach requires huge

manual tuning, also appropriate features needs to be selected and the link classifier has to be created. ACHE framework was projected where crawlers could learn patterns of relevant links automatically and accordingly adapt their focus while crawling the net. This method compact work of manual setup and tuning which was a major drawback of FFC crawlers. Adaptive learning strategy efficiently manages the exploration of acquired knowledge with the exploitation of links with unidentified patterns, this makes crawling process robust and corrects biases introduced in the learning process. Since this crawler learn from scratch, it is able to get harvest rates that is equal or even extra than the manually configured crawlers, so the ACHE framework reduce the effort to configure a crawler manually. ACHE enhanced FFC by incorporating an adaptive link learner and automatic feature selection method. G. Almpandis, C. Kotropoulos and I. Pitas [8] presented a latent semantic indexing classifier that incorporates link analysis along with text content in order to redeem and index the domain explicit net documents. Gunjan H. Agre and Nikita V. Mahajan [9] introduced extraction of links on the basis of search criteria or keyword. It fetches links of web pages that comprises searched keyword in their resources. It gives priority only to those pages and does not fetch web pages unrelated to search keyword. It enhances search effectively with more accuracy, thus providing high optimality as compared to the long-established web crawler. The main characteristic of focused crawler is to collect topic related pages. The earlier crawling experience can help crawler to build knowledge base and learn from it, so that it can get better its performance. Niran et al [10] presented an algorithm that built knowledge bases using seed URLs, keywords and URL guess. These knowledge bases helps crawler to learn and produce the result in more well-organized way. The knowledge bases are incrementally built from the log of previous crawling. Seed URLs allow the sycophant to collect as many related web pages as possible. Keywords support the crawler to recognize proper documents. URL prediction enables the crawler to guess the relevancy of the content of unvisited URLs. Qingyang Xu et al [11] introduced a new general structure for focused crawler that is based on “relational subset discovery”. To describe the significance in between the pages that are not visited in the crawl frontier predicates are used, after that first order classification rules are introduced using subgroup discovery method. Then the learned relational rules guide crawler with sufficient support and confidence.

A new advance for predicting the links that lead to relevant pages based on a Hidden Markov Model (HMM) was presented by Hongyu Liu et al [12]. The method includes three stages: user data collection, user modelling using sequential pattern learning, and focused crawling. Initially internet pages are collected while browsing through web. Then these WebPages are clustered, and then hierarchical linkage pattern within pages from dissimilar clusters is then used to learn sequence of pages that lead to target pages. The Hidden Markov Model (HMM) is used for learning procedure. During the crawling process the precedence of links are decided on basis of a learned estimate of how likely the page will guide to a target page. The performance was compare with Context-Graph crawling during experiments it was establish that this approach performed better than Context-Graph crawling.

III. METHODOLOGY

The hidden web contents are the information content that cannot be indexed by search engines as they stay behind searchable web interfaces. Current hidden web directories mostly have less coverage of related web resources which degrade their ability. A Crawler goes across a no of variety of web pages during a crawling process. Hence to do efficient crawling and wide coverage , ranking and prioritizing links of different sites is necessary. The objective of this method is to extract hidden web information with wide coverage for hidden web resource and uphold efficient crawling for focused web crawler. Spell checker is apply to input for getting correct input and then reverse searching is performed for getting center pages of known website. Incremental site prioritizing algorithm use to precisely rank the relevant links. User is then provided with the top ranked links available in searchable form.

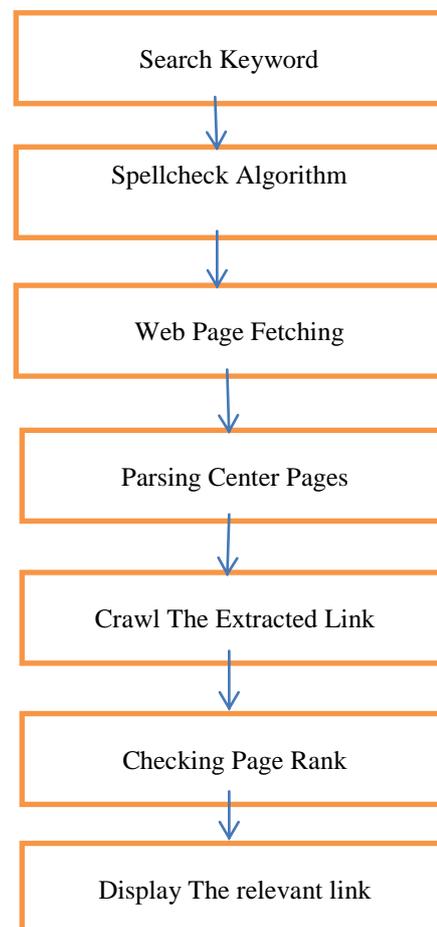


Fig. 1. Flow of Proposed System

The method existing in this paper is based on the framework proposed by Feng Zhao, Jingyu Zhou, Chang Nie, Heqing Huang and Hai Jin [1]. The system consist of a structure with two stages to address the problem of searching for hidden-web resources. The two stages framework of Crawler contains, site locating and in-site exploring as shown in Figure2. The site locating stage discover sites relevant to the given topic. Depend upon the similarity and frequency of sites with known seed sites. Site Similarity means how many content are similar to the known sites in database And The site frequency

determine the number of presences of sites into another sites. The site frequency is defined in equation 1.

$$SF(s) = \sum I_i \quad (1)$$

where $I_i = 1$ if s appears in a known hidden web site, else $I_i = 0$.

If a site is consider as topic relevant, further operation of in-site exploring is done to see searchable forms. An adaptive learning algorithm does online feature selection and build link rankers automatically using those features. Each feature context can be represented as array of terms having specific weight. The weight of term t is represented as vector of term . And weight w of term was given in equation 2.

$$W_{t,d} = 1 + \log t_{t,d} \quad (2)$$

Where $t_{t,d}$ denote the no of count of appearance of term t in document d . Adaptive learning method updates this information collected successfully during crawling.

a) Site Locating - This stage finds the related sites for specified topic by classifying and ranking them. Site Collecting strives to maximize the number of hidden websites and minimize the number of visited URLs. To maximize the number of hidden websites links in the fetched pages is not enough because websites generally contain less links to other sites. As a solution to this problem reverse searching and incremental two-level site prioritizing strategy are used to get additional sites. Reverse searching is performed using seed site in the crawler list (Site collecting).The links are extracted after parsing of the page getting from the search engine. These pages are then downloaded and analyzed to find if the links are related. Then all related links found are intermediate output. To achieve incremental crawling and reach wide coverage Incremental Site Prioritizing is done using High priority Queue H_q and Low priority Queue L_q . The knowledge that is the deep websites or links gained through crawling are recorded. Then, unvisited sites are assigned to Site Frontier and are prioritized by Site Ranker and visited sites are included into the fetched site list. Once Site frontier gets enough sites, Site Ranker allocates a score to the sites depending on how relevant they are to the topic. Site Classifier function will classify the sites according to the search query. If site is found to be related to the topic, crawling process is carried on. Else, that site is abstained and a new site is drawn from the frontier in its place. Topic relevance of a site is based on the matters of its homepage. The homepage content of newly fetched site is collected and parsed by removing stop words and stemming. Then a feature vector is created for the site the resulting vector is given as input to the Classifier to choose if the page is relevant to the topic or not.

b) In-Site Exploring - In-Site Exploring gather the searchable forms in the ranked web sites provided by site locator. Link Ranker offer high relevance score to the links pointing to the pages with searchable form. It gives priorities to the links so that crawler can find searchable forms rapidly. Form Classifier classify the irrelevant and non search able forms

c) Feature selection - Feature selection determines new search patterns. Feature space of hidden web sites (FSS) and Feature space of links (FSL) are used for constructing feature space and training data. The feature space of hidden web sites (FSS) is given as, $FSS = (U, A, T)$

where U, A, T are vectors which maps to the feature context of URL, anchor and text around URL of the hidden web sites. The feature space of links of a site (FSL) is given as,

$FSL = (P, A, T)$

where A and T are anchor and text around links of deep web sites and P is a vector related to the path of the URL , because all the links of a particular site have same domain.

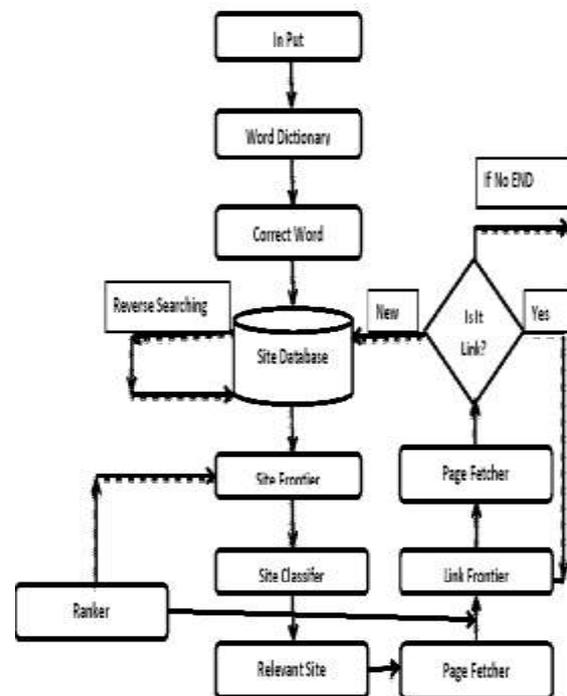


Fig. 2. System Architecture for Hidden(Deep) Web Crawling

A. Algorithms

Following are two algorithms used for deep web crawling. Here reverse searching algorithm proposed by Feng Zhao et al [1] is used as base. Before applying this algorithm spell checker method is used for giving correct input to the system.

a) Reverse searching algorithm:-

Reverse searching algorithm

uses Seed sites and harvested deep web sites (Dws) and Web pages (Wp) to maximize the number of deep web sites in crawling process and provide the relevant sites.

Input : Seed sites and harvested deep websites

Output: Relevant sites

Step 1. While of candidate $W_s < \text{Threshold}$ do

Step 2. $W_s = \text{get}(Dws)$

Step 3. $\text{ResultPage} = \text{ReverseSearch}(W_s)$

Step 4. $L = \text{ExtractLinks}(\text{ResultPage})$

Step 5. For $\forall L$ do

$(W_p) = \text{DownloadPage}(L)$

$\text{Relevant} = \text{Classify}(W_p)$

Step 6. If Relevant then

$\text{RelevantSites} = \text{ExtractUnvisitedSite}(W_p)$

Step 7. Provide Relevant W_s

End

b) Incremental Site Prioritizing algorithm:-

This algorithm uses learned patterns of deep websites to achieve broad coverage on websites W_s and provides the topic

relevant links on the basis of crawling status Cs and page rank PR.

- Step 1. While SiteFrontier ≠ empty do
- Step 2. If Hq ≠ Empty then
- Step 3. Hq.Add All(Lq)
- Step 4. Lq.Clear()
 - end
- Step 5. (Ws) = Hq.poll()
- Step 6. Relevant = ClassifySite(Ws)
- Step 7. If Relevant then PerformInSiteExploring(Ws)
 - Output forms and OutOfSiteLinks
 - SiteRanker.rank(OutOfSiteLinks)
- Step 8. If forms ≠Empty
 - Hq.Add (OutOfSiteLinks)
 - else
- Step 9. Lq.Add(OutOfSiteLinks)
- Step 10. Check(Cs)
- Step 11. Calculate(PR)
- Step 15. Sort(PR,Cs)

Step 16. Provide the relevant Links in descending order.
 Here Fi is input function, Fi= {Q,N,Ws} Where Q=Search keyword by the user , N = Number of search results required , Ws = Web Sites , Fo is output function,
 Fo- {L} Where L= Links Retrieved,
 F= {Pf ,PR,Rr} Where Pf (Ws,L) - Fetch Webpages and new Links.

PR(Ws,L) - Ranking of webpages
 $PR(A)=(1-\delta) + \delta (PR(T1)/C(T1) + \dots + \delta (PR(Tn)/C(Tn))$ (3)
 where,
 PR(K)= PageRank of page K.
 PR(Ti)= PageRank of pages Ti which link to page A.
 C(Ts)= Number of outbound links on page
 δ = damping factor which can be set between 0 and 1.

IV. EXPERIMENTAL RESULTS AND ANALYSIS

The proposed system implementation is done on a machine having 4 GB RAM , Dual core processor, 4 GB RAM and 500GB Hard disk memory for storing URL data is required. Internet connection speed required is minimum 1Mbps for fast hidden web crawling. For security reason it neglects some encrypted or protected URLs from crawling process by hard coding their URL names. Minimum threshold value for number of deep websites is indicated.

To calculate the performance of our crawling framework, we compare Proposed Crawler to the SCDI(site-based crawler for deep web interfaces) .

SCDI : This system is similar to Smart Crawler, SCDI shares the same stopping criteria with Smart Crawler. Only difference is , SCDI follows the out-of-site links of related sites by site classifier without employing incremental site prioritizing technique. It also does not employ reverse searching for collecting sites. prioritizing strategy for sites and links.

Smart Crawler: This is our proposed crawler for harvesting deep web interfaces. Smart Crawler uses an offline-online with offline-online learning to train link classifiers. It leverages learning results for site ranking and link ranking.

Fig 3 shows comparison result between SCDI and proposed system.

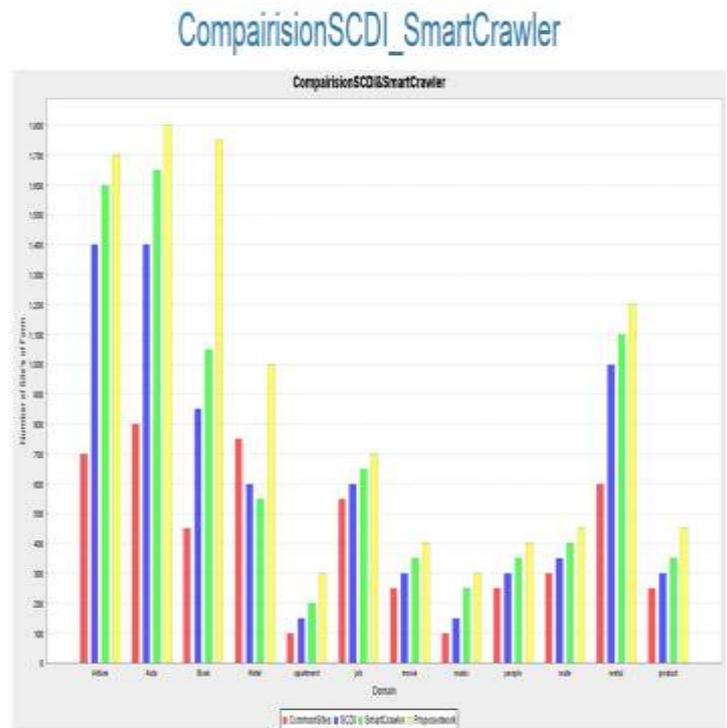


Fig. 3. Comparison of Smart Crawler with SCDI method

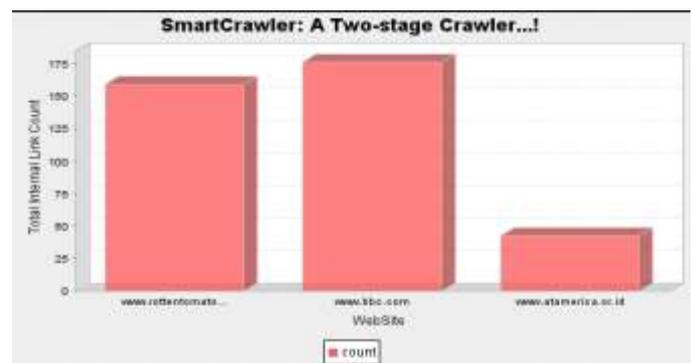


Fig. 4. Total no. of internal sites related to that sites

Table 1
 Experimental result

id	Search key word	ExTime in Mili Sec	Total number of main sites crawled
1	Music	21523902363	25
2	Apartment	60982853685	15
3	Airfare	25043328930	21
4	Book	79943724779	30
5	People	69693107146	21
6	America	17209862449	19

Fig. Using different keywords of variable length experiment was done to explore the efficiency of proposed system. Google API is used to acquire the center web sites. The links around those center websites are crawled using the Incremental Site Prioritizing with the reverse search technique to get more links relevant to the topic. The results gained as per experimental

method explained above are given in Table 1. Number of links crawled shows the number of hidden web links harvested during the crawling process. For key word people crawler crawled no of link and display 21 main sites and for each sites crawler perform crawling and displayed all internal link of that related sites in fig 5. Similar results were found for keywords Music, Apartment, Airfare and Book.

V. CONCLUSION

In this system an enhanced framework is proposed for hidden web harvesting framework for hidden-web interfaces. The method proposed here is expected to provide efficient crawling and more coverage of hidden web. So it is a focused crawler, its searches are topic related and it can rank the fetched sites. Reverse searching and Incremental Site Prioritizing prioritize the relevant links according to the page ranks and displays links with higher page rank, so that it achieves accurate results. The proposed system is expected to achieve more harvest rates than other crawlers

REFERENCES

- [1] Feng Zhao, Jingyu Zhou, Chang Nie, Heqing Huang, Hai Jin, "SmartCrawler: A Two-stage Crawler for Efficiently Harvesting Deep-Web Interfaces", IEEE Transactions on Services Computing 2015.
- [2] YeyeHe , Dong Xin , VenkateshGanti, "Crawling Deep Web Entity Pages", ACM 2013.
- [3] "Information Retrieval", by David A. Grossman and OphirFrieder.
- [4] "Modern Information Retrieval", by Ricardo Baeza-Yates and Berthier Ribeiro-Neto.
- [5] SoumenChakrabarti, Martin Van den Berg, and Byron Dom, "Focused crawling: a new approach to topic-specific web resource Discovery", Computer Networks, 1999.
- [6] L. Barbosa and J. Freire., "Searching for Hidden-Web Databases", In Proceedings of WebDB, pages 1-6, 2005.
- [7] 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.
- [8] G. Almpandis, C. Kotropoulos, I. Pitas, "Combining text and link analysis for focused crawling-An application for vertical search engines", Elsevire Information Systems 2007.
- [9] Gunjan H. Agre, Nikita V. Mahajan, "Keyword Focused Web Crawler", IEEE sponsored ICECS 2015.
- [10] NiranAngkawattanawit and ArnonRungsawang, "Learnable Crawling: An Efficient Approach to Topic-specific Web Resource Discovery", 2002.
- [11] QingyangXu , WanliZuo, "First-order Focused Crawling", International World Wide Web Conferences 2007.
- [12] Hongyu Liu, Jeannette Janssen, EvangelosMilios, "Using HMM to learn user browsing patterns for focused Web crawling", Elsevire Data Knowledge Engineering 2006.