



iJRASET

International Journal For Research in
Applied Science and Engineering Technology



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Volume: 3

Issue: V

Month of publication: May 2015

DOI:

www.ijraset.com

Call: ☎ 08813907089

E-mail ID: ijraset@gmail.com

Focus: Accustom To Crawl Web-Based Forums

M.Nikhil¹, Mrs. A.Phani Sheetal²

¹ Student, Department of Computer Science, GITAM University, Hyderabad.

² Assistant Professor, Department of Computer Science, GITAM University, Hyderabad

Abstract: Due to availability of sufficient data on web, searching has a significant impact. A web crawler is a automatic scheduled program from the huge downloading of web pages from World Wide Web and this process is called Web crawling. To cluster the web pages from World Wide Web a search engine uses web crawler and the web crawler collects this by web crawling. Due to restrictions of network bandwidth, time- swallow and hardware's a Web crawler cannot download all the pages, it is necessary to select the important imperative ones as early as possible during the crawling process and avoid downloading and visiting many irrelevant pages. This paper analysis the researches on web crawling algorithms used on searching.

Index Terms— EIT path, forum crawling, ITF regex, page classification, page type, URL pattern learning.

I. INTRODUCTION

Traditional search engines allow users to submit a query proposing, as output, an or-dered list of pages ranked according to a particular matching algorithm. The under-lying Information Retrieval model's goal is to allow users to find those documents that will best help them meet information needs and make it easier to accomplish their information-seeking activities. The query can be considered the user's textual descrip-tion of the particular information request. If the engine works in the internet domain, a software system usually named *crawler* traverses it, collecting HTML pages or other kinds of resources. It exploits the hyperlink structure in order to find all the destination anchors (or targets) reachable from a given starting set of pages through the outgoing links. General-purpose search engines employ crawlers to collect pages covering different topics. At query time, the engine retrieves subsets of pages that are more likely to be related to the user current needs expressed by means of sets of keywords. Methods of the user needs are usually not employed; therefore the search results are not personalized



Fig. 7.1. Taxonomy of approaches to build specialized search engines, as shown in [80].

for the user. Basically, two users, with different interests, knowledge and preferences, obtain the same results after having submitted the same query. A first step toward a better search tool is developing *specialized search engines*, which provide tailored information on particular topics or categories for focused groups of people or individual users. The heavy constraints of a general-purpose search engine, i.e., indexing billions of pages and processing thousands of queries per second, are no longer required for these kinds of tools. New techniques can be included to represent Web pages and to match these representations against the interests of users, e.g., algorithms based on Natural Language Processing (NLP), usually avoided due to the computational resources needed. There are several approaches to build specialized search engines, as shown in Fig. 7.1. Query modification and re-ranking exploit traditional search tools, filtering their content by augmenting user queries with keywords, or re-ordering the results, removing unwanted resources. Specialized search engines are also based on focused indexes, which contain only the documents related to the given topics of interest. To retrieve and index those documents, it is possible to meta-search specialized databases, or perform an autonomous search at query time. The most interesting technique is to perform focused crawling on the Web. It concerns the development of particular crawlers able to seek out and collect subsets of Web pages that satisfy some specific requirements. In particular, if the goal is to collect pages related to a given topic chosen by the users, the crawlers are usually named *focused* or *topical* [2] (see Fig. 7.2). Focused crawlers are also employed in different domains from specialized IR-based search engines, but usually related to the retrieval and monitoring of useful hypertextual information. The focused crawling approach entails several

International Journal for Research in Applied Science & Engineering Technology (IJRASET)

advantages in comparison with the other approaches employed in specialized search engines. Performing an autonomous search at query time considerably delays the retrieval of result lists. Meta-searching provides results from existing general-purpose indexes that often contain outdated versions

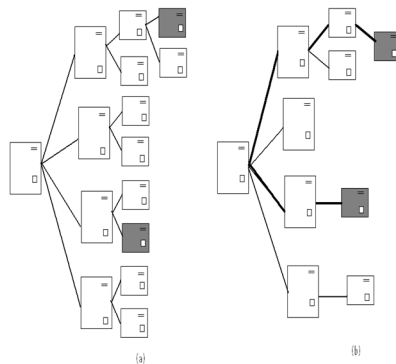


Fig. 7.2. Focused crawling attempts to find out and download only pages related to given topics (b), while standard crawlers employs traditional strategies (a), e.g., breadth-first search or further techniques to balance the network traffic on Web servers.

of the available documents. Due to the reduced storage needs, focused crawling can employ techniques to crawl part of the deep Web (dynamically generated Web pages not accessible by search engines, see Sect. 7.2) or to keep the stored information fresh updating it frequently. If a focused crawler includes learning methods to adapt its behavior to the particular environment and its relationships with the given parameters, ex: the set of retrieved pages and the user-defined topic, the crawler is named *accustom*. Non-*accustom* crawlers are usually based on classifiers whose learning phase ends before the searching process starts. Even though some of them employ hypertextual algorithms, such as HITS - which lead to better results, making more information available - the adaptability is actually not manifest. Sometimes *accustom* crawlers provide a sort of feedback mechanism where useful information is extracted from the analyzed documents, and the internal classifier updated consequently. Other approaches can explicitly model the set of pages around the topical ones. Such models capture crucial features that appear in valuable pages, and describe the content of the pages that are frequently associated with relevant ones. The searching process may benefit from those models, obtaining better overall crawl performance. *Accustom* focused crawlers are key elements in personalizing the human-machine interaction. Conventional non-*accustom* focused crawlers are suitable for communities of users with shared interests and goals that do not change with time. In this case, it is easy to recognize the requested topics and start a focused crawl to retrieve resources from the Web. The *accustom* crawler's advantage is the ability to learn and be responsive to potential alterations of the representations of user requirements. This could occur when users do not know exactly what they are looking for, or if they decide to refine the query during the execution if the results are not deemed interesting. Therefore, *accustom* focused crawlers are more suitable for personalized search systems that include a better model of the information needs, which keeps track of user's interests, goals, preferences, etc.. As a consequence, *accustom* crawlers are usually trained for single users and not for communities of people.

II. RELATED WORK

Vidal et al. proposed a tale approach for learning regular expression patterns of URLs that lead a crawler from an entry page to target pages. Target pages were found through comparing the DOM trees of pages with a preselected trial of target page. It is very frequent but it only works for the particular site from which the sample page is drawn. It is essential to repeat the same process for every time for the new site. However, it is not relevant for large-scale crawling. In contrast, our proposed approach FoCUS which learns URL patterns across multiple sites and automatically finds a forum's entry page given a page from the forum. Guo et al. and Li et al. are comparable to our work. However, they did not mention "how to discover and traverse the URLs". Li et al. developed some heuristic rules to discover URLs. But, the rules are very particular and it can only be applied to specific forums powered by the particular software package in which the heuristics were conceived. Unfortunately, according to the Forum Matrix [3], there is lot of discomparable forum software packages used on the Internet. Wang et al. presented an algorithm to address the traverse path selection problem. They introduced the structure of skeleton link and page-flipping link. Skeleton links are "the most significant links supporting the structure of a forum site." Importance is determined by the informativeness and coverage metrics.

International Journal for Research in Applied Science & Engineering Technology (IJRASET)

Page-flipping links are defined using the connectivity metric. By identifying and only following skeleton links and page-flipping links, they demonstrated that iRobot can achieve effectiveness and coverage. According to our supervision, the sampling strategy and informativeness estimation is not stout and tree-like traversal path is not possible. Traversal path does not tolerate more than one path from a starting page node to a same ending page node. Another related work is in the vicinity of our work which presented to avoid duplicate detection. Forum crawling also desires to remove duplicates. However, this content based duplicate detection does not have competent bandwidth, it can only be carried out when pages have been downloaded. URL-based duplicate detection [7] is not supportive. In forums, index URLs, thread URLs, and page-flipping URLs have specific URL patterns. Thus, in our paper, by learning patterns of index URLs, thread URLs, and page-flipping URLs and adopting a simple URL string de-duplication technique (e.g., a string hashset), FoCUS can be easily avoided duplicates without any duplicate detection. To advance the unnecessary crawling, industry standards such as “no follow” [4], Robots Exclusion Standard (robots.txt) [6], and Sitemap Protocol [5] have been introduced here. By inserting the “rel” attribute with the “no follow” value, page authors can inform a crawler that the destination content is not endorsed. However, it is intended to diminish the effectiveness of search engine spam, but not meant for jamming the access to pages. A proper way is robots.txt. It is designed to specify what pages a crawler is allowed to visit or not. Sitemap [5] is an XML file that lists the URLs along with additional metadata including update time, change frequency and efficiency etc. Generally, the intention of robots.txt and Sitemap is to facilitate to be crawled effectively. So they may be important to forum crawling. However, it is difficult to deal such files for forums as their content continually changes

III. PROPOSED SCHEME

In this we discussed about our proposed scheme and how to instrument, it. In this section, we provide all the method in separate module with detailed description such as synopsis of anticipated scheme, ITF Regexes Learning, Online Crawling and Entry URL Discovery.

A. Overview Proposed Scheme

In this section we present architectural diagram for our anticipated scheme in Fig 3. It consists of two major parts: the learning part and the online crawling part. The learning part first learns ITF regexes of a given forum from automatically constructed URL training examples. The online crawling part then applies learned ITF regexes to crawl all threads efficiently.

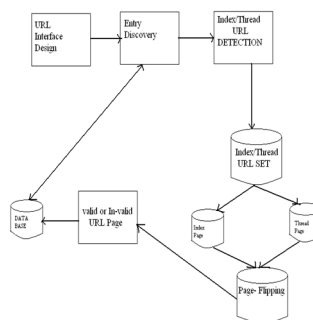


Fig 3 : Architecture

Page Type: In this module, we classified the forum pages into following page types.

Entry Page: The main of a forum, which contains a list of boards and is also the lowest familiar ancestor of all threads.

Index Page: A page of a board in a forum, which mostly contains a table-like structure and which contains information of a board. The list-of board page, list-of-board and the thread page, and the board page are all index pages.

Thread Page: A page of a thread in a forum that contains a list of posts with user generated content belonging to the comparable discussion.

URL Type: In this module, we discuss about types of URL

International Journal for Research in Applied Science & Engineering Technology (IJRASET)

Index URL: A URL that is on an entry page or index page and points to an index page. Its anchor text shows the title of its target board.

Thread URL: A URL that is on an index page and points to the thread page. Its anchor text is title of the destination thread.

B. ITF Regexes Learning

In this section, we learn about ITF regexes, FoCUS which adopts two-step supervised training procedure. The first step is training sets construction. The second step is regexes learning.

1) URL Training Sets: The goal of URL training sets construction is to automatically construct the sets of highly precise index URL, thread URL, and page-flipping URL strings for ITF regexes learning. We use a comparable process to construct index URL and thread URL training sets since they have very comparable properties with the exception of the types of their destination pages.

2) Learning ITF Regexes: In this sub-module, we have shown how to construct index URL, thread URL, and page-flipping URL string training set. We also elucidate how to learn ITF regexes from these training sets. Vidal et al. applied URL string generalization. For example, given URLs as follows (the top four URLs are encouraging while the bottom two URLs are pessimistic):

Instead, we apply the method introduced by Koppula et al. which is advanced to deal with pessimistic examples.

C. Online Crawling

In this module, we perform online crawling using a breadth-first strategy (actually, it is easy to adopt other strategies). FoCUS first pushes the entry URL into a URL queue; next it fetches a URL from the URL queue and finally downloads its page; and then it pushes the outgoing URLs which are coordinated with any learned regex into the URL queue. Accustom FoCUS repeats this step until the URL queue is empty or other conditions are satisfied. FoCUS only needs to apply the learned ITF regexes on innovative outgoing URLs in newly downloaded pages to making the more proficient for online crawling. FoCUS does not need to group outgoing URLs, classify pages, recognize page-flipping URLs, or learn regexes again for that forum.

D. Entry URL Discovery

In this module, an entry URL needs to be precise to start the crawling process. In particular in web-scale crawling, manual forum entry URL bad notation is not practical. Forum entry URL finding is not a trivial task since entry URLs vary from forums to forums. We developed a novel heuristic rule to stumble on entry URL as a baseline. The heuristic baseline tries to stumble on the following keywords ending with “/” in a URL: forum, board, community, bbs, and debate. If a query is found, the path from the URL host to this keyword is extracted as its entry URL; if not, the URL host is retrieved as its entry URL. To make the FoCUS more practical and scalable, we design a simple yet frequent forum entry URL discovery method based on some techniques.

IV. RESULT

The main goal of our experiment is to download relevant forum content from the web with minimal overhead. This can be achieved by using index-thread-page-flipping algorithm that can be used to recognize index, thread, or page-flipping URLs. Some of the screen shots of our experiment are

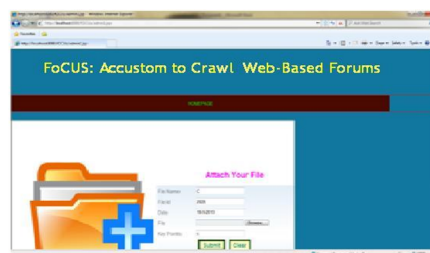


Fig 51 . File Upload Page

International Journal for Research in Applied Science & Engineering Technology (IJRASET)



File Id	File Name	Uploaded On	Key Points
7021	Object	1-8-2013	staticBehavior
9906	Scope	1-8-2013	object-oriented programming
7913	oops	1-8-2013	Object,Class,Data,abstraction,Data Encapsulation,Interface,Polymorphism,Dynamic Binding
6703	Class	1-8-2013	objects
9726	Data Abstraction	1-8-2013	programming language

Fig 5.2 File Details:



Fig 5.3 User Home Page

In this module, we report the results of our implemented paper. FoCUS achieved 99 percent precision and 99 percent recall. The low standard deviation also designates that it is not sensitive to sample pages.

V. CONCLUSION

In the paper we present the method crawler which downloads and stores web pages, frequently for a web search engine. The fast growth of web poses more challenges to search for suitable link. We also symbolize the technique of FOCUS which are developed to extract only the relevant web pages of interested topic from the Internet. The design of FOCUS is capable to evaluate the text which found on a link with the input text file. The crawler uses pattern recognition and generates the number of times the input text exists in the text establish on a link. Particular attention has been given to accustom focused crawlers, where learning methods are able to adapt the system behavior to a particular environment and input parameters during the search. Evaluation results show how the whole searching process may benefit from those techniques, enhancing the crawling performance. Adaptivity is a must if search systems are to be personalized according to user needs, in particular if such needs change during the human-computer interaction.

REFERENCES

- [1]. Jingtian Jiang, Xinying Song "FoCUS: Learning to Crawl Web Forums" IEEE TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING, VOL. 25, NO. 6, JUNE 2013.
- [2]. Chakrabarti, S.: Recent results in automatic web resource discovery. ACM Computing Surveys **31**(4es) (1999) 17
- [3]. R. Cai, J.-M. Yang, W. Lai, Y. Wang, and L. Zhang, "iRobot: An Intelligent Crawler for Web Forums," Proc. 17th Int'l Conf. World Wide Web, pp. 447-456, 2008.
- [4]. Y. Guo, K. Li, K. Zhang, and G. Zhang, "Board Forum Crawling: A Web Crawling Method for Web Forum," Proc. IEEE/WIC/ACM Int'l Conf. Web Intelligence, pp. 475-478, 2006.
- [5]. K. Li, X.Q. Cheng, Y. Guo, and K. Zhang, "Crawling Dynamic Web Pages in WWW Forums," Computer Eng., vol. 33, no. 6, pp. 80-82, 2007.
- [6]. G.S. Manku, A. Jain, and A.D. Sarma, "Detecting Near-Duplicates for Web Crawling," Proc. 16th Int'l Conf. World Wide Web, pp. 141-150, 2007.
- [7]. M.L.A. Vidal, A.S. Silva, E.S. Moura, and J.M.B. Cavalcanti, "Structure-Driven Crawler Generation by Example," Proc. 29th Ann. Int'l ACM SIGIR Conf. Research and Development in Information Retrieval, pp. 292-299, 2006.



10.22214/IJRASET



45.98



IMPACT FACTOR:
7.129



IMPACT FACTOR:
7.429



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24*7 Support on Whatsapp)