

Information Retrieval based on Content and Location Ontology for Search Engine (CLOSE)

Niranjan Kumar¹ and S G Raghavendra Prasad²

¹ Rashtreeya Vidyalaya College of Engineering

Received: 16 December 2013 Accepted: 1 January 2014 Published: 15 January 2014

Abstract

This paper mainly focuses on the personalization of the search engine based on data mining technique, such that user preferences are taken into consideration. Clickthrough data is applied on the user profile to mine the user preferences in order to extract the features to know in which users are really interested. The basic idea behind the concept is to construct the content and location ontology's, where content represent the previous search records of the user and location refer to current location of user. SpyNB is the approach used to mining the user preferences from the Clickthrough data. The ranked support vector machine (RVSM) is performed on the searched results in order to display results according to user preferences by considering Clickthrough data.

Index terms— SpyNB, personalization, ontology, RSVM, non-geographic search, geographic search, search engine optimization (SEO), personalized information retrieval

In the modern information retrieval system, the results that are found should be more accurate to query submitted by the user, and also efficiency should be considered.

In order to solve the problems that are faced by the current search engine technology such as retrieving results that are irrelevant to the search query, the order in which they are displayed should be considered. According to Hele-Mai Haav [1] to solve problem of information retrieval in current information retrieval systems it should be improved by intelligence to manage the effective retrieval, filtering and presenting relevant information. So two main information retrieval models are classified as, keyword based information retrieval model and concept based information retrieval model. The indexing terms and Boolean logical queries are used in keyword based model, where indexing may be automatic or manual, when Boolean query are taken into consideration the frequency of occurrence is taken into account.

Context-aware system [2], depending on the user's relevancy the information/services is provided. For instance consider the keyword apple, it can mean as a fruit or it can mean as a mobile and laptops by Apple Company. When the query is submitted by two different users, irrespective of their interest same results are displayed for both users, if one user is interested only on apple accessories, for him both relevant and irrelevant information are displayed in random order. The information for what the user is looking may be in same document else somewhere in the overall document. The current system performs word to word matching of the search query.

Another instance in search engine is searching for places based on current location of the user. For example, if the user current location is Jaynagar and user trying to search restaurant near by current location, the search engine must show the restaurant which are near to the current location of the users and rest of the restaurant location other than jaynagar should be given next preference. The detailed discussion related to geographic and non-geographic search is given in proposed system section.

The main aspects that should be considered in information retrieval system is to reduce the complexity involved in query execution [3] such that performing lexical analysis, stemming process on the user query and construction of index terms. This paper focuses on search engine optimization (SEO) by reducing the complexity in the user query execution.

The rest of the paper is organized as: -In section II literature survey is carried out by surveying previous paper present, such that what are the technologies currently used to optimize the search engine. In section III technique to reduce the complexity for optimization of search query are studied. In section IV detailed view of implementation. In section V experimental evaluation and in IV Conclusion and future enhancements are discussed.

1 II.

2 Literature Survey

M. Rami Ghoran [4] studied that for every query that is submitted by the user he will get the relevant and irrelevant information for that query. So they classify the personalized information retrieval (PIR) system into three scopes: Individualized system, community-based system and aggregate-level system.

When individualized system is considered the system adaptive [6] decision are taken such that, the user interest and preferences are taken into account while Performing the search operations, while this approach leads to true to true personalization but it has some drawback such as:

Fresh start, when user is new to system his/her interest should be tracked and some time user may not compromise to share personal information with the system. Community-based system [7] describes sharing of the information among several users/models. The data enrichment technique such as clustering technique is used in grouping of the similarity among various users. Using some similarity criteria the users among the web can be grouped into one model, so that results for this community can be personalized. Aggregate-level system [8] where information gathered is represented in the form of summary for purpose of analysis. The common parameters such as age are considered to form clusters. For example a site selling music CD's may advertise certain CD's based on the age of the users and data aggregate for their age group. Online analytic processing (OLAP) is the simple type of data aggregation.

Browser also provides certain level of personalization by storing the cookies and recently visited web hyperlinks in the buffers. When the user is in static place browser will provide certain level of personalization, but when user place changes dynamically buffer contents are no more used.

For this purpose the new technique can be taken into consideration, such that each user's interest is maintained in the server buffer so that where ever user requests some result in form of query this can be compared with user interest buffer and relevant information can be retrieved from the system by minimizing unrelated results. When the user is new to system and enters any query for the first time the preferences for location is taken along with search keyword and search operation is performed. The keyword of the query is searched in the server and relevant results are fetched and displayed as the results. When the user clicks on some links, Click through data will be recorded. Later when the user searches for the same keyword, the previously visited pages will be displayed first with higher ranked pages and, if there is are any new links they will be ranked in lower order.

3 III.

4 System Design

Spy NB [9] is the algorithm used to fetch the user Click through data, and these are transformed to vectors for further process. The Ranked support Vector machine (RSVM) training is performed on the vectors for Re-ranking of search results according to user preferences. The detailed description about Spy NB and RSVM is given in implementation part.

The system mainly concentrates on building the method of ontology for all the possible keywords. The word can have different meaning in different context [2].

For example when the keyword "JAVA" is considered, in several perspectives it mean as the programming language, but by the name JAVA there is an island in Indonesia, and java coffee is referred to as a coffee beans.

When the two users submit the query both will get similar results either list of Java Island or list of java Where $s f(c_i)$ is the web snippet frequency of the keyword/phrase in the query Q , n is the total number of web snippet and $|c_i|$ is the number of terms in the keyword/phrase c_i . If the support of the keyword/phrase c_i is higher than threshold \hat{I}^T (where threshold \hat{I}^T is set by user), then we consider c_i as the concept for query Q .

In this system the value of \hat{I}^T is set to 5 because, if \hat{I}^T value is assigned with lesser value than for each search, ranking should be updated this leads to consume more time for reordering of links. assigned with larger value than perfect personalization cannot be achieved.

The following two prepositions are adopted to find relationship between concepts for ontology:

? Similarity: The two concepts which coexist more in the search results can be considered or represented as the same topic of interest. If occurrence of document $c_i, c_j > \hat{I}^T$ (where \hat{I}^T is the threshold) then c_i and c_j can be considered as similar.

? Parent-Child Relationship: specific concepts appear with general terms, but backtracking is not true. If the preference of c_i and $c_j > \hat{I}^T$ then we can conclude that c_i is child of c_j . possible concept space determined for the keyword/phrase "Nokia" while Click through data will determine the preferences on the concept based.

The concept space for the query "nokia" consists of different types of models such as E-series, N-Lumina etc. When E-series is taken into consideration both has similarity that they belong to same parent.

Content space for the query "Nokia" consists of "N1100", "E-series", "6600", and so on. If the user is interested in E-series and clicks on the page containing price, the Click through of the links are captured. These Click through data is considered as the positive preferences and vector is constructed.

When the same query is issued by the same user later the vector is transferred to server by transforming this content vector into content weight vector to rank the search result according to user preferences.

Location Ontology: The approach of the location ontology [13] [14] [15] is quite different from the construction of content ontology. Following assumptions are made i.e., the parent-child relationship cannot be accurately derived for the location ontology. To construct the vector [15] Content Ontology: The concept works on extracting the keywords/phrase from the web snippets by eliminating all the stems in the query Q. The content ontology is classified differently to different users based on their interest. The co-existence of the keyword in the query Q is calculated to find similarity among the user interest by using following support and confidence rule [3]:

Clickthrough data: It is the process of recording the links or advertisement that is clicked by the user(s), for the purpose of determining which link is viewed how many times. The system makes use of these Clickthrough [10] data in personalizing each specific user's interest by maintaining the records for each user in the database. In formal language it can be defined as, it is triplets of (Q, R, C) where Q is the query, R is the ranking order in which it is displayed and C is the set of URLs that are clicked by the users. To achieve personalization the system is classified into two distinct levels namely, content ontology and location ontology [12]. The detailed descriptions about two levels are elaborated in below section:

Bangalore, "Jaynagar/Bangalore/Karnataka/India", is associated with the document d. The construction of the vector for the location ontology is similar to that of the content ontology. The Clickthrough data is transferred to the server and transformed as the location vector and this vector is used to rank the user preferences.

IV.

5 Implementation

In this section technique that are used to personalize the search engine are discussed in detail. First, when the query q is entered by the user, look for previous records if previous search results are found then apply Content ontology concept else if the user is new then accept the query q and apply Location ontology concept.

Ranking algorithm will rank the results according to the user preferences by calculating the weight of both content and location concepts, for keyword/ key phrase. The content weight of all posts for particular keyword is considered in calculating the ranking order.

The vector support machine is constructed for training the user preferences, loop is entered when the ranking operation is started, and the number of count is recorded for the link whenever the user clicks on it. When the post reaches the minimum threshold value then it will gain a higher order value as compared from rest of the post. The formal representation for performing these is depicted below:

6 Return Result

Spy Naive Bayes (SpyNB) algorithm is used to collect the Clickthrough data. This algorithm will maintain two sets called positive set Ps and negative set Ns. Where $P = \{\text{Links that are clicked by the users}\}$ $N = \{\text{Links not clicked by the users}\}$ Algorithm 1: CLOSE (U_i, q, L) // Input: User identity U_i , Query q and Current location of User L. // Output: Results for query with user preferences. 1. Accept the Query q from user where $q \in \{A-Z, a-z, 0-9\}$ 2. Filter the post (documents) using the keyword q If ($\text{Post}(d_i) == \text{compare}(q)$) 3. If (check user profile U_i for previous records) Next algorithm will be related to searching keyword based on Content ontology. Algorithm 2: Content-Ontology (U_i, q) // Input: User -Identity, and corresponding Query q. // Output: Return Results to CLOSE Next algorithm will be related to searching keyword based on Location ontology. Algorithm 3: Location-Ontology (U_i, q, L) Algorithm 4: SpyNB(s) // Input: Post matched for Query q. // Output: Feature vector for Post 1. Compare S with the user record.

7 Experimental Evaluation

The Table 1 gives the dataset of the content ontology construction for some of the keywords. The table mainly consists of unique code for particular root keyword, name of keyword and parent of the corresponding keyword [17]. In the experimental evaluation "Hotel" is the root word and it has four children such as "Reservation", "Facilities", "Meeting Room", and "Party hall", similarly for others also constructed.

Similarly Table ?? gives the dataset of the location ontology construction for some of the locations.

The table mainly consists of location code, Location name, latitude, longitude and parent of location. When location is considered, boundary value of 11 values is taken into consideration.

8 Table 2 : Statistic of Location Ontology

In posting of documents the related information are stored by entering the root and location for which it belongs. In this case Hotel "comfort" comes under Bangalore city for which India will be root, and so on others are posted.

When user enters the query q , the searching process will be carried out as mentioned in the implementation section by invoking several techniques. When the corresponding documents are found, and previous records of users are analyzed, the ranking support vector machine is performed on the posts that are matched by the keyword or query q .

Table 3 gives the RSVM calculation for the Keyword "jaguar" for two different users, it can be observe from the table that two user have their own preferences in choosing the link.

Later, when two users search for same keyword then threshold value changes and ranking of their search results will be altered.

9 Global

10 Conclusion and Future Enhancement

We can conclude that the CLOSE system will provide better search results as compared to rest of the search engines by considering the users Content and location concepts. CLOSE system will take user preferences in minimizing the possible time for retrieving search results. RSVM training will be performed for each individual user profile, so that system will come to know in what the user is really interested.

As a future enhancement it can be extended by considering time as one of the parameter to even more optimize the search results. The sessions can also be considered as one of the parameter, so that when user stop work at particular instance, later when user get into system, at moment where user stopped working or viewing content of some documents, from that session it should be started (with respect to two or more different systems).

11 VII.

12 Global



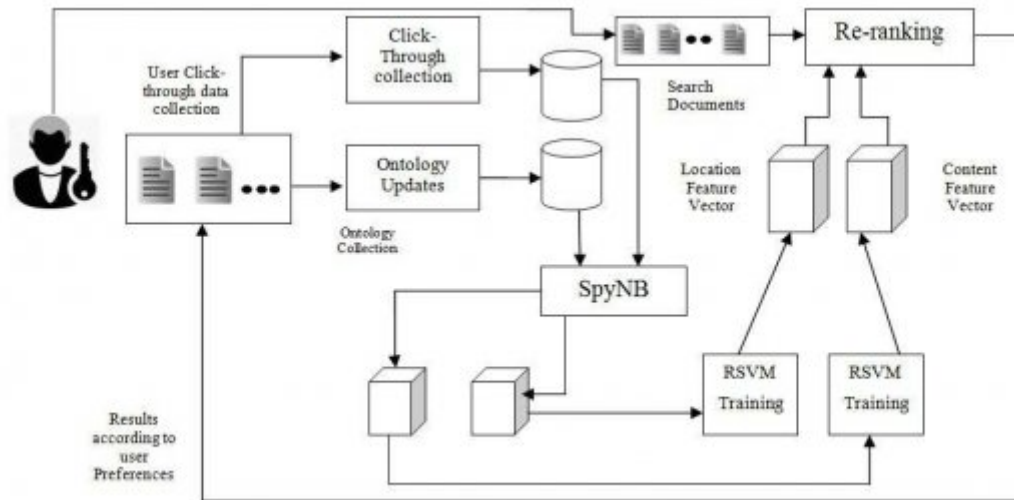
Figure 1: Figure 1 :

1 2 3

¹© 2014 Global Journals Inc. (US) Information Retrieval Based on Content and Location Ontology for Search Engine (CLOSE) coffee beans is displayed or list of java programming is displayed, but one user expecting only about island and other only programming language. The system mainly

²© 2014 Global Journals Inc. (US) Information Retrieval Based on Content and Location Ontology for Search Engine (CLOSE)

³© 2014 Global Journals Inc. (US)



1

Figure 2: Fig 1

N 1100 Features	q=Nokia	6600 E 5631	Level 0
	E Series		N Lumina
	E 5630		Level 1
	Price		Level 2
	Parent-child relationship		
	Similarity		

Figure 3:

1

Unique Code Keywords Parent 101 Hotel 0 102
 Reservation 101 103 Facilities 101 104 Meeting Room
 103 105 Party Hall 103 106 Animal 0 107 Jaguar 106
 108 Lion 106 109 Car 0 110 Jaguar 109 111 BMW 109
 112 Black Jaguar 107 113 Elephant 106

Figure 4: Table 1 :

Algorithm 5: RSVM (count, post_code)

// Input: count for each click is taken as the input.

// Output: Ranking order of the posts.

1. For i 0 to total_post-1 do
2. Content_weight_count count.
3. Calculate the Content weight for particular keyword. P_code Post_code

4. Content_weight (%)

5. Final_content_weight

6. P1

7. P2 P1-
100

8. location_weight_parameter 9. Final_rank Final_content_weight + location_weight_parameter

Locat

Code

1 1

12 1

123 1

124 1

1231 .

1232 1

13 7

1

2 1

21 1

a

1

1

22 1

23 1

Unique KeywordsParent

Code

101 Hotel 0

102 Reservation 101

103 Facilities 101

104 Meeting 103

Room

105 Party 103

Hall

106 Animal 0

107 Jaguar 106

108 Lion 106

109 Car 0

110 Jaguar 109

111 BMW 109

112 Black 107

Jaguar

113 Elephant 106

Figure 5:

3

Figure 6: Table 3 :

.1 Acknowledgement

Foremost, I would like to express my sincere gratitude to my guide Mr. S G Raghavendra Prasad Assistant Professor, ISE Dept, RVCE, for the continuous support of my M. Tech study, for his patience, motivation, enthusiasm, and immense knowledge. His guidance helped me in all the time of writing this technical paper.

Besides my guide, I would like to thank the rest of my M.Tech committee: Dr. Jitendranath Mungara PG Dean, ISE Dept, RVCE, and Dr. Cauvery N K. HOD ISE Dept, RVCE.

.2 Author

[IEEE RIVF ()] , *IEEE RIVF* 2013. p. .

[Papagelis and Zaroliagis ()] ‘A Collaborative Decentralized Approach to Web Total=11 Total=866 Search’. Athanasios Papagelis , Christos Zaroliagis . *IEEE Transaction on Systems, Man, and Cybernetics* 2012. 42 (5) p. .

[Haav and Lubi] *A Survey of Concept based Information Retrieval Tools on the Web*, Hele-Mai Haav , Tanel-Lauri Lubi . (White paper)

[Bouramoul ()] ‘An ontology-based approach for semantics ranking of the web search engines results’. Abdelkrim Bouramoul . *IEEE* 2012. p. .

[Ghorab] *Centre for Next Generation Localisation Knowledge & Data Engineering Group*, M Rami Ghorab . p. . (Personalised information retrieval: survey and classification)

[Bhatia ()] ‘Context-aware Personalized Mobile Web Search Techniques-A Review’. Deepika Bhatia . *IJCSIT International Journal of Computer Science and Information Technologies* 2011. 2 (5) p. .

[Jain and Mahajan ()] ‘Data Mining Based on Semantic Similarity to Mine New Association Rules’. Sandeep Jain , Aakanksha Mahajan . *Global Journal of Computer Science and Technology Software & Data Engineering* 2012. 12. (Issue 12 Version 1.2)

[Mobasher ()] ‘Discovery and Evaluation of Aggregate Usage Profiles for Web Personalization’. Bamshad Mobasher . *ACM Transaction on Data Mining and Knowledge Discovery* 2002. 6 (1) p. .

[Veningston and Shanmugalakshmi ()] ‘Enhancing personalized web search re-ranking algorithm by incorporating user profile’. R K Veningston , Shanmugalakshmi . *IEEE* 2012. p. .

[Sharji] ‘Enhancing the Degree of Personalization through Vector Space Model and Profile Ontology’. Al Sharji , Safiya . *IEEE Computing and Communication Technologies* p. 2013. (Research. Innovation, and Vision for the Future (RIVF))

[Sun ()] ‘FoSSicker: A Personalized Search Engine by Location-Awareness’. Mingyang Sun . *IEEE* 2012. p. .

[Mishra ()] ‘Improving Mobile Search through Location Based Context and Personalization’. Varun Mishra . *IEEE* 2012. p. .

[Ng ()] ‘Mining User Preference Using Spy Voting for Search Engine Personalization’. Wilfred Ng . *ACM Transaction on Internet Technologies* 2007. 7 (3) p. .

[Baeza-Yates and Ribeiro-Neto ()] *Modern Information Retrieval: The Concepts and Technology behind Search*, R Baeza-Yates , B Ribeiro-Neto . 2013. (Pearson Edition)

[Li Qing-Shan ()] ‘Ontology based User Personalization Mechanism in Meta Search Engine’. Li Qing-Shan . *IEEE* 2012. p. .

[Shen ()] ‘Query Enrichment for Web-query Classification’. Dou Shen . *ACM Transactions on Information Systems* 2006. 24 (3) p. .

[Goel ()] ‘Search Engine Evaluation Based on Page Level Keywords’. Shikha Goel . *IEEE* 2013. p. .

[Raval and Kumar ()] ‘SEReLeC (Search Engine Result Refinement and Classification) -A Meta Search Engine based on Combinatorial Search and Search Keyword based Link Classification’. Vishwas Raval , Padam Kumar . *IEEE* 2012. p. .

[Arora and Kant ()] ‘Techniques for Adaptive websites and Web Personalization without any user effort’. Kanika Arora , Kamal Kant . *IEEE Students conference on Electrical, Electronics and Computer Science* 2012.