

Classified Average Precision Approach for Inferring User Search Goals with Feedback Sessions



Shaik Samsher
M. Tech Student,
Department of CSE,
Guntur Engineering College, Yanamadala.



N. Bhagya Lakshmi, M.Tech
Assistant Professor,
Department of CSE
Guntur Engineering College, Yanamadala.

Abstract:

Ambiguous queries are submitted in search engine by different users with different search goals. The analysis of user click through logs can be useful in finding the precise search results. The user click through logs contains the information about the user search information. By analyzing the user click through logs the feedback sessions are constructed.

The pseudo documents are generated by representing the feedback sessions for clustering. A novel approach for user and query dependent feedback sessions for user search results. The CAP is formulated to evaluate the performance of user search goal inference. This can be very useful in improving search engine efficiency. We are reducing risk factor value in proposed approach.

Keywords:

User Seek Objectives, Input Sessions, Pseudo-Reports, Cap, Vap,

INTRODUCTION:

Web mining is one of the applications of data mining techniques to discover knowledge from the web. In web search, users are submitted queries to the search engines to get relevant information. But many search engines results are not informative and failed to produce results according to the user search goals. Users are usually giving some vague keywords representing their interests in their minds. Such keywords do not match with the results produced by the search engines.

Many works about user search goals analysis should be carried out. Some users give ambiguous queries to the search engines (e.g. Apple, jaguar, the sun etc.) they get mostly the irrelevant results. User search goals are classified as Navigational and Informational, the queries that seek a single website or webpage and queries that reflect the intent of the user to perform a particular transaction respectively. Many related works have been carried out according to the web search applications and the user search goals. In previous works, clustering is done on a set of top ranked results.

The user search logs information is not analyzed and the feedback sessions are not considered. Analyzing the clicked URLs only from the web search logs. They only identify whether a pair of queries belong to the same goal or mission and does not care about what the goal is in detail.

Semantic based web search for a particular query and the similarity between the words are carried out. Various algorithms such as star clustering algorithm, k-means clustering algorithm are used for clustering the pseudo documents but it also does not cluster the relevant information according to the user search goals. In clustering the cluster labels discovered are also not informative.

User search goal is the information on different aspects of a query that users want to obtain. Information need is a user's wish/desire to obtain the relevant information to satisfy his need. To cluster web search results, the URLs are analyzed by extracting the titles and snippets. But all those works produced noisy results and does not obtain the user search goals precisely.

When more irrelevant and relevant results are produced by the search engines it is time consuming to browse. In this paper, the user submits the query into the browser. The search engine searches the relevant information according to the user query. The user actions are stored in the user click through logs. From the user click through logs each and every session is analyzed and generates the feedback session. The feedback session contains both the clicked and unclicked URLs and the last clicked URL

in a single session. The feedback session contains the URLs and the click sequence. By analyzing the feedback sessions, the pseudo documents are generated. The pseudo documents contains the keywords that are most clicked in a session. Likewise the pseudo documents are clustered using the clustering algorithm. The user search goals are obtained according to the feedback sessions. The restructure result is produced for the user query based on the user search goal. The CAP evaluation can be done for each user search goal and the clustering can be done to find the optimal number of users.

RELATED WORK:

Recently, numerous works have been carried out to infer the so-called client objectives or aims of a question [13]. Yet in actuality, their works have a place with question arrangement. A few works examine the query items returned by the web crawler straightforwardly to adventure diverse inquiry angles [6]. Be that as it may, question viewpoints without client criticism have constraints to enhance web crawler significance. A few works take client criticism into record and investigate the diverse clicked URLs of an inquiry in client navigate logs straightforwardly, in any case the quantity of diverse clicked URLs of a inquiry may be not huge enough to get perfect results.

Wang and Zhai grouped questions and scholarly parts of these comparative inquiries, which tackles the issue to some degree. Notwithstanding, their strategy does not work on the off chance that we attempt to find client look objectives of one single inquiry in the question bunch as opposed to a group of comparative questions. For instance, in [12], the inquiry “auto” is grouped with some dissimilar questions, for instance, “auto rental,” “utilize auto,” “fender bender,” and “auto sound.”

Thus, the distinctive parts of the inquiry “auto” have the capacity be learned through their structure. nevertheless, the inquiry “utilize auto” in the group can likewise have distinctive angles, which are hard to be adapted by their system. Some different works present seek objectives and missions to recognize session limit progressively [11]. In any case, their strategy just recognizes whether a couple of inquiries fit in with the same objective or mission what’s more does not give a second thought what the objective is in point of interest. A former use of client navigate logs is to get client implied input to expand preparing information when learning positioning capacities in data recovery.

Thorsten Joachims did numerous chips away at how to utilize certain criticism to enhance the recovery quality [8], [9], [10]. In our work, we consider input sessions as client certain criticism and propose a novel improvement technique to join both clicked and unclicked URLs in input sessions to figure out what clients truly require and what they couldn’t care less. One application of client inquiry objectives is rebuilding web indexed lists. There are likewise some related works centering on arranging the list items [6], [8]. In this paper, we induce client seek objectives from client navigate logs and rebuild the query items as indicated by the construed client look objectives.

PROPOSED FRAMEWORK:

Our framework comprises of two sections separated by the dashed line. In the upper part, all the feedback sessions of a query are initially extricated from user click-through logs and mapped to pseudo-documents. At that point, user seek objectives are gathered by clustering these pseudo-documents and delineated with some essential words. Since we don’t have a clue about the careful number of user look objectives ahead of time, a few distinctive qualities are attempted furthermore the ideal quality will be dictated by the feedback from the bottom part.

In the bottom part, the first query items are rebuilt focused around the user look objectives derived from the upper part. At that point, we assess the execution of rebuilding query items by our proposed assessment rule CAP. Also the assessment result will be utilized as the feedback to choose the ideal number of user hunt objectives in the upper part.

Queries are submitted to search engines to represent the information needs of users. Ambiguous queries contain one or several polysemous terms. Query ambiguity is one of the main reasons for poor retrieval results (difficult queries are often ambiguous). User Click-through data log contains data about the interactions between users and Web search engines.

It is one of the most extensive (yet indirect) surveys of user experience. The user search information's are stored in the user click through logs. From the user click through logs the feedback sessions are constructed. The proposed feedback session consists of both clicked and unclicked URLs and ends with the last URL that was clicked in a single session.

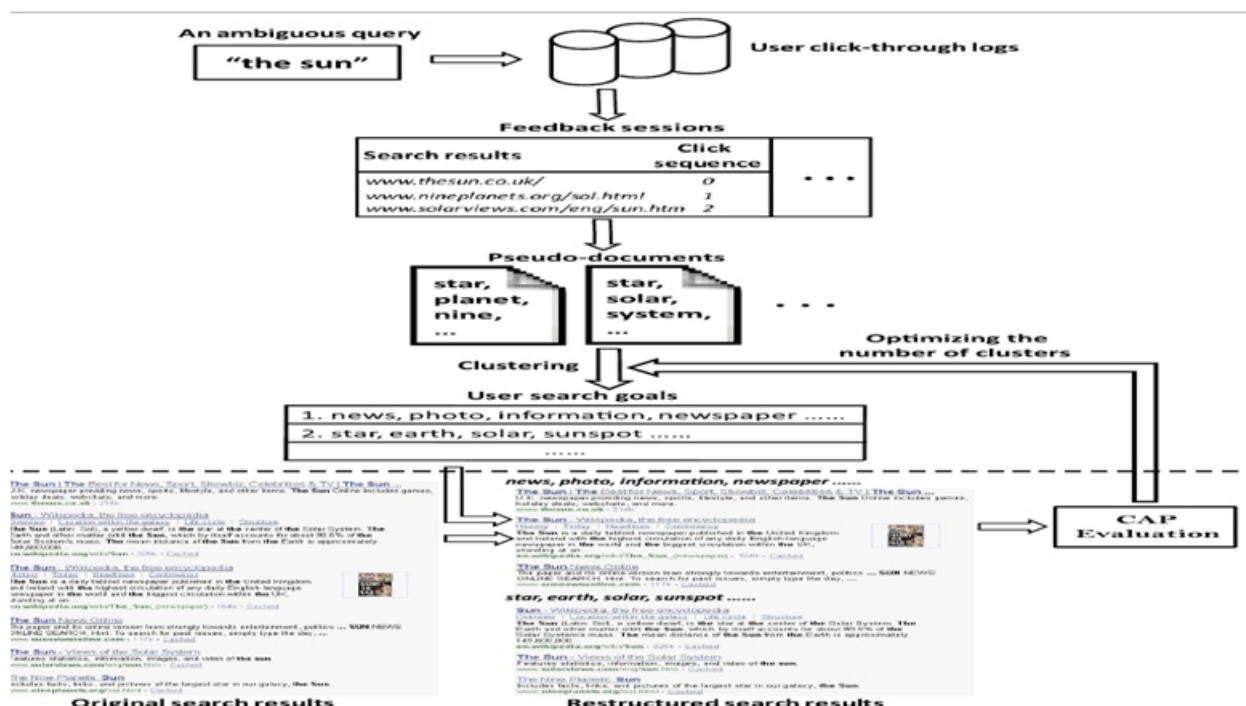
The feedback sessions is based on a single session, although it can be extended to the whole session. The feedback session contains the URLs with the click sequence. A novel optimization method to combine the enriched URLs in a feedback session to form a pseudo-document. This can effectively reflect the information need of a user.

The URLs are enriched from the feedback sessions based on the click sequence. The enriched URLs with more value in click sequence are mapped to pseudo-documents.

The pseudo documents are depicted with some keywords based on the URL. At last, cluster these pseudo documents to infer user search goals and depict them with some keywords.

ANALYZING USER CLICK THROUGH LOGS:

The user click through logs is analyzed for each session to propose a feedback sessions. The feedback session is the better representation for the user click through logs. It is more efficient than analyzing the user click through logs directly. For a single query each and every session is analyzed and represents the feedback session. The feedback session is based on a single session although it can be extended to the whole session. An ambiguous query is that it gives more than one meaning. So the precise results according to the user search goal are difficult to obtain.



Feedback Session:

The first process in reorganizing the search result is the feedback session representation. Feedback session consist the list of URLs up to the URL that was clicked by the user at last in a single session. All the unclicked URLs before the last clicked URL in a single session is also included because those URLs also has been browsed and analyzed by the user. Therefore, these unclicked URLs must also be included for the feedback. From this feedback session, the clicked URLs represent what information the user entail and the unclicked URLs reflect what information the user do not require. The URLs that are present after the last clicked URL cannot be taken as a part of feedback because it is not certain whether the user have scanned those URLs or not. Feedback session cannot be used directly for user search goal inference because it varies from that of the user click-through logs. So, it should be represented in some other forms in order to infer the user search goals efficiently.

It can be represented in various forms. Binary vector representation is one of the popular ways of representing the feedback session. It consists of 0's and 1's where "0" represents the unclicked URL and "1" represents the clicked URL in a single session. This method cannot be used when more feedback sessions are considered because diverse feedback sessions may have unusual aspects. The vague keywords can be used to represent the user interests for a query. But these keywords cannot be used for representing the feedback session because they are usually hidden and not expressed clearly. Therefore, pseudo-documents can be used to infer the goals of the user. The feedback sessions are mapped to the pseudo-documents. These documents can be formed by enriching those URLs present in the feedback session. Enriching the URLs can be done by adding the title and a short snippet in a small text paragraph for the same URLs.

Pseudo Documents:

The pseudo documents are not the legitimate documents. The URLs in the feedback sessions are enriched by some format. The URLs are formatted by removing the stop words and the stemming words. It is the icon of showing the information about the whole document by some keywords.

The documents are created by the number of occurrences of the keywords. The keywords which are having the more frequency are grouped together. The pseudo documents contain the keywords that are retrieved from the URLs in the feedback sessions. Using the Meta tag information the URLs are enriched. The Meta tag contains the most important keywords about the entire document information.

EVALUATIONS OF RESTRUCTURE SEARCH RESULTS:

The evaluation of user search goals can be done using the CAP (CLASSIFIED AVERAGE PRECISION). The classified average precision is the calculation of precision of documents. Because from user click-through logs, we can get implicit relevance feedbacks, namely "clicked" means relevant and "unclicked" means irrelevant.

A possible evaluation criterion is the average precision (AP) which evaluates according to user implicit feedbacks. AP is the average of precisions computed at the point of each relevant document in the ranked sequence. VAP is the voted average precision which can be used for grouping the dissimilar documents for the particular user query search. Risk is the mapping of similar and dissimilar documents for the particular user query. If there is a similarity then the mapping value is 0 and if there is no similarity between VAP and risk then the mapping value is 1.

AVERAGE PRECISION:

A possible evaluation criterion is the average precision (AP) which evaluates according to user implicit feedbacks. AP is the average of precisions which is computed at the point of each relevant document in the ranked sequence, shown in

$$AP = \frac{1}{N^*} \sum_{r=1}^{N^*} rel(r) \frac{R_r}{r}$$

Where, N is the number of relevant (or clicked) documents in the retrieved ones, r is the rank, N is the total number of retrieved documents, rel(r) is a binary function on the relevance of a given rank, and R_r is the number of relevant retrieved documents of rank r or less.

RISK:

It is the AP of the class including more clicks There should be a risk to avoid classifying search results into too many classes by error. So we propose the Risk.

$$Risk = \frac{\sum_{i,j=1(i < j)}^m d_{ij}}{C_m^2}$$

Voted AP(VAP):

VAP of the modernized search result the AP of class 1, It is defined as

$$VAP = \frac{1}{NC} \sum_{r=1}^{NC} rel(r) \frac{R_r}{r}$$

where N is the total numeral of retrieved documents with class label one , rel() is a binary function on the relevance of a given rank, and R_r is the number of relevant retrieved documents of rank r or less Classified Average Precision (CAP) Extend VAP by introducing the above Risk and propose a new criterion .

Classified AP(CAP):

$$CAP = VAP * (1 - risk)^y$$

Where r is used to adjust the influence of Risk on CAP. CAP select the AP of the class with the aim of user is interested with the most clicks/votes and takes the risk of wrong classification into account.

User Search Goals:

Cluster pseudo-documents by FCM bunch that is easy and effective. Since we have a tendency to don't understand the precise variety of user search goals for every question, we have a tendency to set variety of clusters to be 5 totally different values and perform bunch supported these 5 values, severally.

Once bunch all the pseudo-documents, every cluster will be thought of mutually user search goal. The middle purpose of a cluster is computed because the average of the vectors of all the pseudo-documents within the cluster.

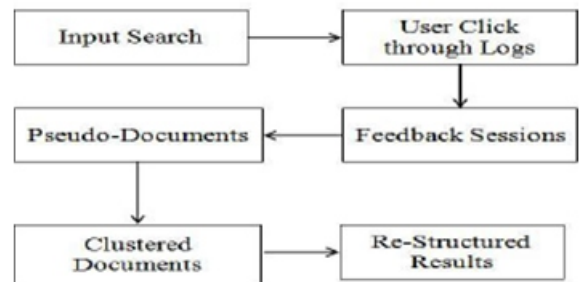


Fig:-overall data flow of the proposed system
Algorithm :-

Input: User s Search Query (Q)

Output: Re-Structured Results (RR)

BEGIN

Get query „Q from user

Populate Q History as a query history dataset with the retrieved

results from the database – DBQueries

If

QHistory is empty

Return Actual Web Search Results „RW

Else

For each Query Instance Q_i in QHistory

Generate Feedback Session FSi

Convert Feedback sessions to Pseudo Documents PDi

End For

End If

Set Web Search Results „RW as

results from web search engine

Set Output „RR using the cluster_function(PDi, Q_i, Q)

Return Re-Structured Results – „RR

END

$$CAP = VAP * (1 - Risk) \quad (1)$$

A single session	click sequence	rel(r)R _r /r
www.sun.co.uk	0	0
www.nineplanets.org	1	1/2
www.solarviews.com	2	2/3
en.wikipedia.org	0	0
www.thesunmagazine.org	0	0
www.space.com	0	0
en.wikipedia.org/the sun (newspaper)	3	3/7
www.nasa.gov	0	0
www.nasa.gov/worldbook	4	4/9
AP = 1/4[1/2+2/3+3/7+4/9]=0.510		
Class 1	click sequence	rel(r)R _r /r
www.nineplanets.org	1	1/1
www.solarviews.com	2	2/3
en.wikipedia.org	0	0
www.space.com	0	0
www.nasa.gov	0	0
www.nasa.gov/worldbook	4	3/6
VAP = 1/3[1/1+2/2+3/6]=0.833		
Class 2	click sequence	rel(r)R _r /r
www.sun.co.uk	0	0
www.thesunmagazine.org	0	0
en.wikipedia.org/the sun (new paper)	3	1/3

Object Evaluation and Comparison:

Method I clusters the top 100 search results to infer user search goals [6], [20]. First, we program to automatically submit the queries to the search engine again and crawl the top 100 search results including their titles and snippets for each query.

Then, each search result is mapped to a feature vector according to (1) and (2). Finally, we cluster these 100 search results of a query to infer user search goals by K-means clustering and select the optimal K based on CAP criterion.

Method II clusters different clicked URLs directly [18]. In user click-through logs, a query has a lot of different single sessions; however, the different clicked URLs may be few. First, we select these different clicked URLs for a query from user click through logs and enrich them with their titles and snippets as we do in our method.

Then, each clicked URL is mapped to a feature vector according to (1) and (2). Finally, we cluster these different clicked URLs directly to infer user search goals as we do in our method and Method I.

CONCLUSIONS AND FUTURE WORK:

In this paper, a novel approach has been proposed to infer user search goals for a query by clustering its feedback sessions represented by pseudo-documents. First, we introduce feedback sessions to be analyzed to infer user search goals rather than search results or clicked URLs. Both the clicked URLs and the unclicked ones before the last click are considered as user implicit feedbacks and taken into account to construct feedback sessions. Therefore, feedback sessions can reflect user information needs more efficiently. Second, we map feedback sessions to pseudodocuments

to approximate goal texts in user minds. The pseudo-documents can enrich the URLs with additional textual contents including the titles and snippets. Based on these pseudo-documents, user search goals can then be discovered and depicted with some keywords.

Finally, a new criterion CAP is formulated to evaluate the performance of user search goal inference. Experimental results on user click-through logs from a commercial search engine demonstrate the effectiveness of our proposed methods.

The complexity of our approach is low and our approach can be used in reality easily. For each query, the running time depends on the number of feedback sessions. However, the dimension of Ffs in (3) and (5) is not very high.

Therefore, the running time is usually short. In reality, our approach can discover user search goals for some popular queries offline at first. Then, when users submit one of the queries, the search engine can return the results that are categorized into different groups according to user search goals online. Thus, users can find what they want conveniently.

Future work will be to collaborate query classification and search result combination so that user will get more classified results. Here we are increasing query link length, and we are reducing risk factor value in proposed approach

REFERENCES:

- [1] Zheng Lu, HongyuanZha, Xiaokang Yang, Weiyao Lin, and ZhaohuiZheng, 2013 "A New Algorithm for Inferring User Search Goals with Feedback Sessions" Published by the IEEE Computer Society, pp. 502-522.
- [2] H. Cao, D. Jiang, J. Pei, Q. He, Z. Liao, E. Chen, and H. Li, 2008 "Context-Aware Query Suggestion by Mining Click-Through," Proc.14th ACM SIGKDD Int'l Conf. Knowledge Discovery and Data Mining (SIGKDD '08), pp.875-883.
- [3] X. Li, Y.-Y Wang, and A. Acero, 2008 "Learning Query Intent from Regularized Click Graphs," Proc. 31st Ann. Int'l ACM SIGIR Conf. Research and Development in Information Retrieval (SIGIR '08), pp. 339-346.
- [4] D. Shen, J. Sun, Q. Yang, and Z. Chen, 2006 "Building Bridges for Web Query Classification," Proc. 29th Ann. Int'l ACM SIGIR Conf. Research and Development in Information Retrieval (SIGIR '06), pp. 131-138.
- [5] X. Wang and C.-X Zhai, 2007 "Learn from Web Search Logs to Organize Search Results," Proc. 30th Ann. Int'l ACM SIGIR Conf. Research and Development in Information Retrieval (SIGIR '07), pp. 87-94.
- [6] Giovanna Castellano, A. Maria Fanelli, Corrado Mencar and M. Alessandra Torsello 2007 "Similarity-based Fuzzy clustering for user profiling" Published by the IEEE Computer Society.
- [7] Lazzerini, B. Marcelloni, F.; Cococcioni, M. "A system based on hierarchical fuzzy clustering for web users profiling" Published by the IEEE Computer Society, Print ISBN: 0-7803-7952-7.
- [8] T. Joachims, "Evaluating Retrieval Performance Using Clickthrough Data," Text Mining, J. Franke, G. Nakhaeizadeh, and I. Renz, eds., pp. 79-96, Physica/ Springer Verlag, 2003.
- [9] T. Joachims, "Optimizing Search Engines Using Clickthrough Data," Proc. Eighth ACM SIGKDD Int'l Conf. Knowledge Discovery and Data Mining (SIGKDD '02), pp. 133-142, 2002.
- [10] T. Joachims, L. Granka, B. Pang, H. Hembrooke, and G. Gay, "Accurately Interpreting Clickthrough Data as Implicit Feedback," Proc. 28th Ann. Int'l ACM SIGIR Conf. Research and Development in Information Retrieval (SIGIR '05), pp. 154-161, 2005.
- [11] R. Jones and K.L. Klinkner, "Beyond the Session Timeout: Automatic Hierarchical Segmentation of Search Topics in QueryLogs," Proc. 17th ACM Conf. Information and Knowledge Management (CIKM '08), pp. 699-708, 2008.
- [12] R. Jones, B. Rey, O. Madani, and W. Greiner, "Generating Query Substitutions," Proc. 15th Int'l Conf. World Wide Web (WWW '06), pp. 387-396, 2006.
- [13] U. Lee, Z. Liu, and J. Cho, "Automatic Identification of User Goals in Web Search," Proc. 14th Int'l Conf. World Wide Web (WWW '05), pp. 391-400, 2005.
- [14] X. Li, Y.-Y Wang, and A. Acero, "Learning Query Intent from Regularized Click Graphs," Proc. 31st Ann. Int'l ACM SIGIR Conf. Research and Development in Information Retrieval (SIGIR '08), pp. 339-346, 2008.
- [15] M. Pasca and B.-V Durme, "What You Seek Is what You Get: Extraction of Class Attributes from Query Logs," Proc. 20th Int'l Joint Conf. Artificial Intelligence (IJCAI '07), pp. 2832-2837, 2007.

- [16] B. Poblete and B.-Y. Ricardo, "Query-Sets: Using Implicit Feedback and Query Patterns to Organize Web Documents," Proc. 17th Int'l Conf. World Wide Web (WWW '08), pp. 41-50, 2008.
- [17] D. Shen, J. Sun, Q. Yang, and Z. Chen, "Building Bridges for Web Query Classification," Proc. 29th Ann. Int'l ACM SIGIR Conf. Research and Development in Information Retrieval (SIGIR '06), pp. 131-138, 2006.
- [18] X. Wang and C.-X. Zhai, "Learn from Web Search Logs to Organize Search Results," Proc. 30th Ann. Int'l ACM SIGIR Conf. Research and Development in Information Retrieval (SIGIR '07), pp. 87-94, 2007.
- [19] J.-R. Wen, J.-Y. Nie, and H.-J. Zhang, "Clustering User Queries of a Search Engine," Proc. Tenth Int'l Conf. World Wide Web (WWW '01), pp. 162-168, 2001.
- [20] H.-J. Zeng, Q.-C. He, Z. Chen, W.-Y. Ma, and J. Ma, "Learning to Cluster Web Search Results," Proc. 27th Ann. Int'l ACM SIGIR Conf. Research and Development in Information Retrieval (SIGIR '04), pp. 210-217, 2004.
- [21] D. Beeferman and A. Berger, "Agglomerative Clustering of a Search Engine Query Log," Proc. Sixth ACM SIGKDD Int'l Conf. Knowledge Discovery and Data Mining, pp. 407-416, 2000.
- [22] H. Chen and S. Dumais, "Bringing Order to the Web: Automatically Categorizing Search Results," Proc. SIGCHI Conf. Human Factors in Computing Systems, pp. 145-152, 2000.
- [23] R. Jones and K.L., Klinkner "Beyond the Session Timeout: Automatic Hierarchical Segmentation of Search Topics in Query Logs," Proc. 17th ACM Conf. Information and Knowledge Management, pp. 699- 708, 2008.
- [24] U. Lee, Z. Liu and J. Cho, "Automatic Identification of User Goal sin Web Search," Proc. 14th Int'l Conf. World Wide Web, pp. 391-400, 2005.
- [25] X. Li, Y.Y. Wang, and A. Acero, "Learning Query Intent from Regularized Click Graphs," Proc. 31st Ann. Int'l ACM SIGIR Conf. Research and Development in Information Retrieval, pp. 339-346, 2008.
- [26] B. Poblete and B.Y. Ricardo B, "Query-Sets: Using Implicit Feedback and Query Patterns to Organize Web Documents," Proc. 17th Int'l Conf. World Wide Web, pp. 41-50, 2008.
- [27] X. Wang and C.X. Zhai, "Learn from Web Search Logs to Organize Search Results," Proc. 30th Ann. Int'l ACM SIGIR Conf. Research and Development in Information Retrieval pp. 87-94, 2007.
- [28] B.R. Yates, C Hurtado, and M Mendoza, "Query Recommendation Using Query Logs in Search Engines," Proc. Int'l Conf. Current Trends in Database Technology, pp. 588-596, 2004.
- [29] H.J. Zeng, Z. Chen, W.Y. Ma and J. Ma, "Learning to Cluster Web Search Results," Proc. 27th Ann. Int'l ACM SIGIR Conf. Research and Development in Information Retrieval pp. 210-217, 2004.
- [30] L. Zheng, Z. Hongyuan, Y. Xiaokang, L. Weiyao and Z. Zhaohui, "A New Algorithm for Inferring User Search Goals with Feedback Sessions," IEEE Transactions on Knowledge and Data Engineering, Vol. 25, No. 3, 2013.

Author Profile:

shaik samsheer

obtained the B.Tech degree in Information Technology(IT) from RVR & JC College of Engineering,Chowdavaram,Guntur District.at. Present i am persuing the M.Tech in Computer Science and Engineering(CSE) Department at Guntur Engineering College,Yanamadala,Guntur District.

N.Bhagya lakshmi

is working as an Assistant Professor in the Department of Computer Science and Engineering at Guntur Engineering College,Yanamadala,Guntur District,A.P,INDIA.