

# International Journal of Innovative Research in Science, Engineering and Technology

(An ISO 3297: 2007 Certified Organization)

Vol. 5, Issue 6, June 2016

## Inferring User Search Goals Using Click Through Log

Amruta Sangavikar<sup>1</sup>, Prof. Deepak Uplaonkar<sup>2</sup>

P.G. Student, Department of Computer Engineering, JSPM Narhe, Pune, Maharashtra, India<sup>1</sup>

Assistant Professor, Department of Computer Engineering, JSPM Narhe, Pune, Maharashtra, India<sup>2</sup>

**ABSTRACT:** Identifying or inferring user's search goal from given query is a difficult job as search engines allow users to specify queries simply as a list of keywords which may refer to broad topics, to technical terminology, or even to proper nouns that can be used to guide the search process to the relevant collection of documents. Information needs of users are represented by queries submitted to search engines and different users have different search goals for a broad topic. Sometimes queries may not exactly represent the user's information needs due to the use of short queries with ambiguous terms.

Hence to get the best results it is necessary to capture different user search goals. These user goals are nothing but information on different aspects of a query that different users want to obtain. The judgment and analysis of user search goals can be improved by the relevant result obtained from search engine and user's feedback.

Here, feedback sessions are used to discover different user search goals based on series of both clicked and unclicked URL's. The pseudo-documents are generated to better represent feedback sessions which can reflect the information need of user. With this the original search results are restructured and to evaluate the performance of restructured search results, classified average precision (CAP) is used. This evaluation is used as feedback to select the optimal user search goals.

**KEYWORDS:** User's Goals, Feedback Session, CAP, Evaluations, Information Retrieval.

### I. INTRODUCTION

Various intentions exist in different users mind while searching things over internet. But by using similar keywords for semantically different concepts, internet provides ambiguous data which is vastly available in database. Search engines accept the user keywords as search queries and fetch the entire data relevant to the entered query from the internet database by web crawling and hence has merely no intelligence for entered query's semantic aspect. Major problem with the web crawled data is that the results are not specific to intended goals. For example the query "the bat" will retrieve the data relevant to the cricket bat, the fling bat and the batman also. Hence it is therefore essential to grab different user search goals. Most of the websites over the internet intend to provide users with search results for entered queries. Many a times user queries are ambiguous due to its shorter length i.e. mostly two to three words which increase the ambiguity of the results to be obtained. Such obtained results do not exactly meet the user goals for searching the keyword. Also different search results are obtained from different search engines at different times. The irrelevant results obtained over internet search, then waste users' time by unnecessarily surfing the data over irrelevant data obtained. Thus the proposed system for inferring the user intentions by analyzing the user click logs. The proposed system tries to signify various user goals as clusters. This will efficiently help users to differentiate between relevant and irrelevant data obtained after searching. The user intentions are based on or classified over user keywords entered for obtaining the results. And further based on the clustering results, the results are restructured [3] [5]. To increase the search efficiency over internet, various techniques were invented and proposed by authors like recognition of search results, classification of query, and session limit detection. Below figure depicts the possible scenario of different user intentions for different times.

# International Journal of Innovative Research in Science, Engineering and Technology

(An ISO 3297: 2007 Certified Organization)

Vol. 5, Issue 6, June 2016



Figure 1.(Goal Text)Different user intentions in their mind

These keywords are order less. Keywords here are represented as search goals or goal text. Thus there is no essential way to classify the queries to search engines, and its prevailing since long time that the query formulation has always been a bottleneck for search engines. Most of the document classification is only centered by researchers, which train the various machine learning algorithms using adequately large number of terms. The errand of arranging web questions is diverse in that web inquiries are short, giving not very many inborn elements [7]. Therefore, many approaches make essential use of documents to classify the results obtained by the entered query. For an instance, the user enters the query 'the bat' to the search engine. Generally the obtained results should be relevant to the bat as a cricket playing instrument. But it displays the data related to the bat movie and the bat bird. Although the user gets the expected results but the results are unstructured and hence uses waste a lot of time in searching for relevant results from this unstructured results. On the off chance that client needed to submit inquiry 'the bat' it will firstly demonstrates the consequence of mailing server rather than the bat film.



Figure.2 Different Result for Query

# International Journal of Innovative Research in Science, Engineering and Technology

(An ISO 3297: 2007 Certified Organization)

Vol. 5, Issue 6, June 2016

## II. RELATED WORK

Application of query classification before result retrieval is done in [13]. Query classification was performed before any dataal document was gathered. A simple pre retrieval process is carried out for query classification. Researchers have proposed three distinct techniques for classification. For such methods training the database with the training datasets is required. [14] Agglomerative query clustering emphasizes on mining of the user data depending on user search logs of history and URLs for determining similar contents . To achieve this, user click through log is used as a feedback for obtaining the user goal by maintaining and analyzing the URL click sequence. The major aim of the proposed system is to collect the data about the clicked through data and thereby cluster the clicked as well as unclicked URLs into different clusters as per relevance and thereby compute its similarity. The second paper is question recommendation utilizing mining navigate data. The real go for question recommendation is to enhance the web search tool execution and internet searcher productivity. In any case, the real downside of the question proposal methods is that these systems don't take a shot at connection of the outcomes or URLs acquired. The drawback of this system is that this system only checks whether the query belongs to same context as that of the search goal. A histogram of the search results is created in the Zealous algorithm [12] and the values of results below to the specific threshold are removed while the values of results above the threshold are used in the search goal inference. It eradicates the terms whose values are lower than defined threshold. Generally, search engines accepts the query and checks for the history in search engine log and accordingly gives the search results.

### A. Privacy preserving algorithm

Author centers on collection of queries by making use of users query search logs. Most Frequent item are displayed in the ZEALOUS [12]. User privacy is maintained by using the zealous algorithm. This security is safeguarded regarding clicked log, inquiry and objectives of the clients. Energetic calculation works in two stages, in the primary stage passionate calculation computes the histogram and in second stage ardent kills the things from histogram which have the qualities lower than the edge.

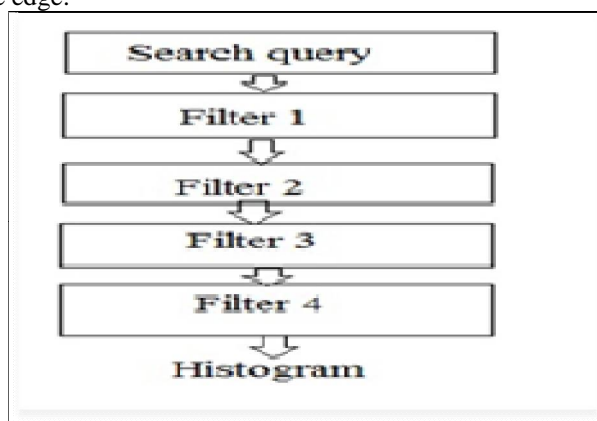


Figure: 3 Flow of Zealous algorithm

Major drawback of the zealous algorithm is that it does not consider user feedback sessions. So the obtained results are much noisy. This system examines two major issues by combining the pre and post retrieval results or classification. So the proposed framework plans and executes a novel technique to recognize the client seek objectives and appropriately group the got unstructured URLs. The proposed strategy is depicted in the accompanying areas. On the implementation of the discrete Fourier transform in the encrypted domain

# International Journal of Innovative Research in Science, Engineering and Technology

(An ISO 3297: 2007 Certified Organization)

Vol. 5, Issue 6, June 2016

## III. PROPOSED SYSTEM

The proposed system designs a new technique to determine the user search goals by logging the user click through sequence. The system makes use of k means clustering algorithm to cluster the results and considers three clusters for the implementation of the same. The system aims in providing the clustered output of the obtained results. And remove the data which are not required or noisy from the searched results. In this system the user enter the query and submits it into the browser. According to the user query the relevant data search by the engine searches. The actions of user are saved in the user click through logs. Each and every session is analysed and generates the feedback session from the user click through logs. According to the feedback sessions there obtain the user search goals. Based on the user search goal the restructure result is produced for the user query. The same query search by the every user with various intensions.

Case in point if client A and B both composed same question in a web look apparatus i. e. web index. Accept their inquiry is 'bat'. The client needs the data about bat and client B needs the data about winged creature Bat. By then as demonstrated by their navigate logs and their looking conduct, the bunching is finished. This bunching have effect in looking for when both clients A and B needs to find same inquiry with different intensions. This navigate log is just the info review of all the eventual outcome of inquiry inquiries. This diagram will help client to find the critical result. Dependent upon this input the pseudo records are made. After that depending on the clients intrigue the navigate report is delivered. Using this report grouping of the client question yield is done. By then applying Cap development method the portrayed yield is appeared. This described yield is just the typical result which client needs to look for.

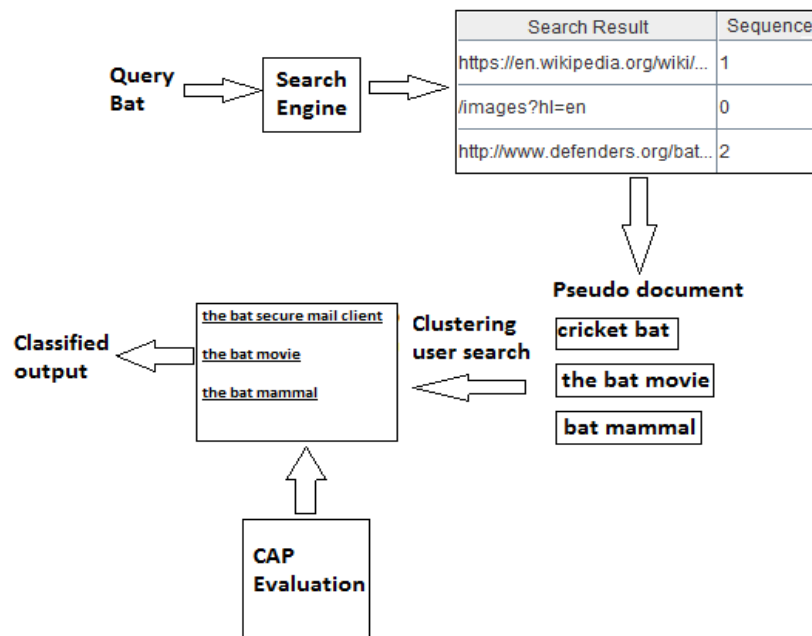


Figure 4. System Architecture

### A. Feedback session:

This is the main strategy for the proposed framework. In this all feedback of the user is get and this qualities are put away in the database in the arrangement of 0s and 1s. This is expected to cluster the URLs for future use. The input sessions are just the clicked and unclicked URL's by the client in the result set. The clicked URLs addresses what clients need and the unclicked URLs identifies with what clients don't require about. The unclicked URLs after the last clicked URL should not be joined into the data sessions since it is not certain whether they were inspected or not. The

# International Journal of Innovative Research in Science, Engineering and Technology

(An ISO 3297: 2007 Certified Organization)

Vol. 5, Issue 6, June 2016

data session can tell what client need is and what kind of result he envision. The data sessions are numbered on the reason of client snap data. The snap plan is escape into session. Right when next time client seeks same question then the session will execute the same navigate game plan to find an accurate result which client needs. An input session is identified with by a little entry of substance that contains its title and some short data. By then, some scholarly systems, for instance, changing each one of the letters to lowercases, stemming and clearing stop words are completed to those substance segments. By then every URL is identified with by some term repeat. By then the weight of each URL is got by some scientific operations. By then these inquiry repeat and URL weight is use to convey pseudo reports.

Search Result	Sequence	Binary
<a href="https://en.wikipedia.org/wiki/">https://en.wikipedia.org/wiki/...</a>	1	1
<a href="/images?hl=en">/images?hl=en</a>	0	0
<a href="http://www.defenders.org/bat...">http://www.defenders.org/bat...</a>	2	1
<a href="http://animals.sandiegozoo...">http://animals.sandiegozoo...</a>	5	1
<a href="http://www.bat.com/">http://www.bat.com/</a>	3	1
<a href="http://www.batcon.org/">http://www.batcon.org/</a>	4	1
<a href="https://www.batcon.org/adopt">https://www.batcon.org/adopt</a>	0	0
<a href="http://www.bats.org.uk/">http://www.bats.org.uk/</a>	6	1
<a href="http://www.bats.org.uk/page...">http://www.bats.org.uk/page...</a>	0	0
<table style="width:auto" bor...	0	0

Figure 5. Clicked Sequence

## B. K Means Clustering Algorithm

In the proposed framework we have utilized K\_means clustering algorithm to cluster the user query output. The k-means algorithm deals with the clickthrough log of the user clicked clustering. The working of the k\_means calculation is as follows.

Let  $X = \{x_1, x_2, x_3, \dots, x_n\}$  be the arrangement of data focuses and  $V = \{v_1, v_2, \dots, v_c\}$  be the arrangement of focuses.

- 1) Randomly select "c" bunch focuses.
  - 2) Calculate the separation between each data point and bunch focuses.
  - 3) Assign the data point to the bunch focus whose separation from the group focus is in particular the group focuses.
  - 4) Recalculate the new bunch focus
  - 5) Recalculate the separation between each data point and new obtained group focuses.
  - 6) If no data point was reassigned then stop, for the most part rehash from step 3).
- Resulting to making after the above strides we will get the grouped yield of the web URLs.

## C. Pseudo Documents

The clicked grouping is used to frame the pseudo archive. Pseudo reports contains the development URLs containing same substance. The effective input session spoke to by pseudo reports. Client may tapped on such a different of connections, so that there may be the making of various criticism sessions. In that all criticism sessions the reports which are having more efficiency than others are known as pseudo archives. In this the snap bunching is re-situated depending on the client clicks. For different rundown things assorted input sessions are kept up. For this I have used one vector known as double vector. The parallel vectors speaks to the technique require for input sessions. With the help of pseudo archive we can without quite a bit of a stretch make derive about client's goals.

For the period of pseudo archives we joins both clicked URL and unclicked URL. By then after the figuring of archive repeat and URL weight the unmistakable match of client's typical result is evaluated. This result is then secured in pseudo record for further future hypothesizing of client need. At whatever point in future client enters same or vital inquiry in web searcher then these pseudo record will convey the result which client needs.

# International Journal of Innovative Research in Science, Engineering and Technology

(An ISO 3297: 2007 Certified Organization)

Vol. 5, Issue 6, June 2016

## D. CAP (Classified Average Precision)

Portrayed ordinary precision is used to evaluate the client list things. This novel system is profitable to choose the best bunch amongst the amount of groups. This will keeps up the metric of client question things. This will chooses client look for goals are determined suitably or not. Depend on upon the criteria used as a part of the CAP we furthermore find the best bunch. In the top we are getting information from the client clicked, clicked infers appropriate and unclicked suggests immaterial. This will offer us to choose some help with being client getting his objective arranged result or not.

## IV. EXPERIMENTAL RESULTS

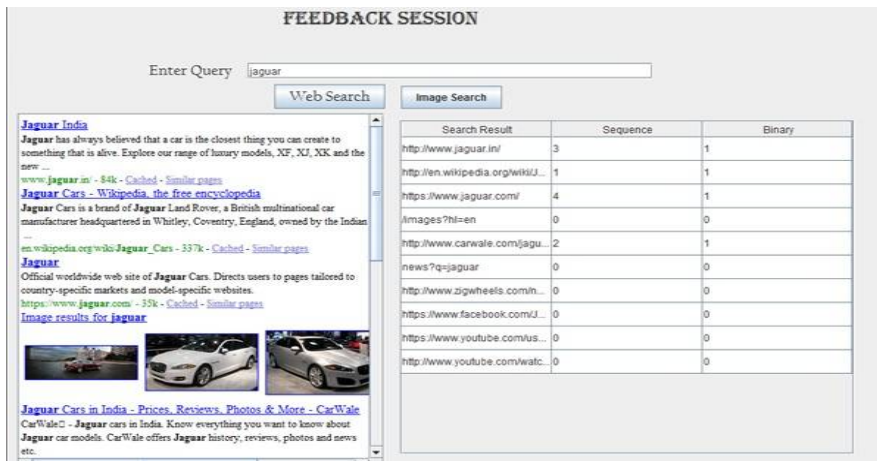


Figure 7. User Query for Feedback

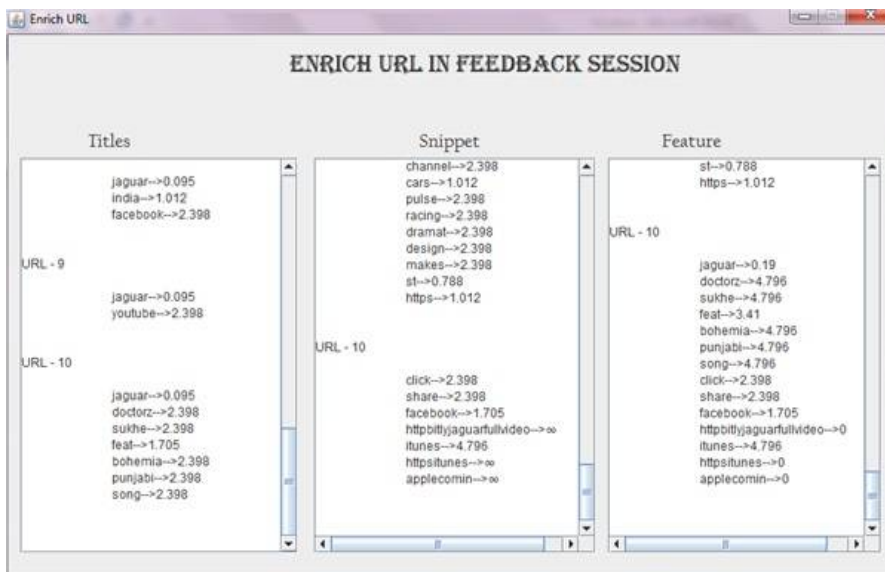


Figure 8. Generating TF-IDF Values of each term

# International Journal of Innovative Research in Science, Engineering and Technology

(An ISO 3297: 2007 Certified Organization)

Vol. 5, Issue 6, June 2016



Figure 8. Generating Pseudo Documents



Figure 9. Reordered Clustered Results

# International Journal of Innovative Research in Science, Engineering and Technology

(An ISO 3297: 2007 Certified Organization)

Vol. 5, Issue 6, June 2016

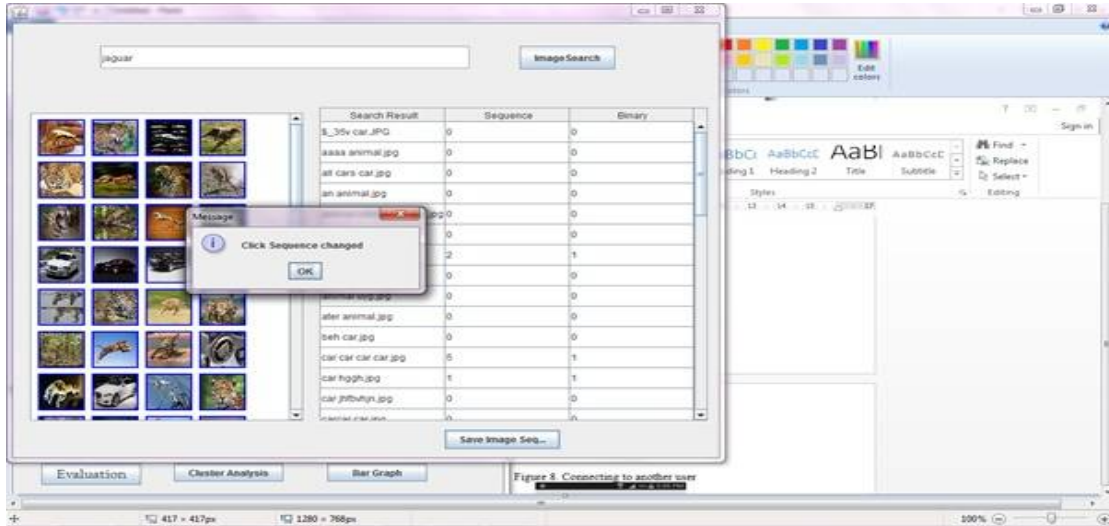


Figure 11. Image Search Feedback

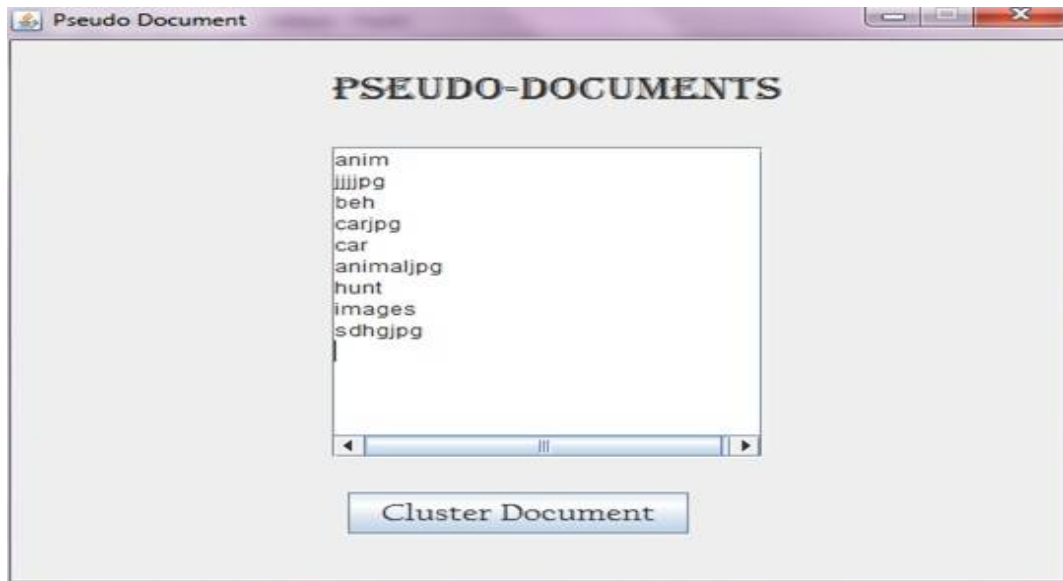


Figure 12. Image search Pseudo Documents.

## V. EXPERIMENTAL ANALYSIS

At the time of experimental analysis of the proposed procedure we will seek different inquiries over the business internet searcher and same question will be fire on our proposed strategy. We will do the examination of the old web searcher inquiry and our systems result. We will do the examination of the result by differentiating particular pursuit calculations. We have added new CAP method to evaluate the execution of the web indexes. The Experimental analysis also provides us with the graphical analysis of the number of terms in which cluster and number of URLs in which Cluster. The figures 14 and 15 represent the graphical view of the clustered terms and urls.



# International Journal of Innovative Research in Science, Engineering and Technology

(An ISO 3297: 2007 Certified Organization)

Vol. 5, Issue 6, June 2016

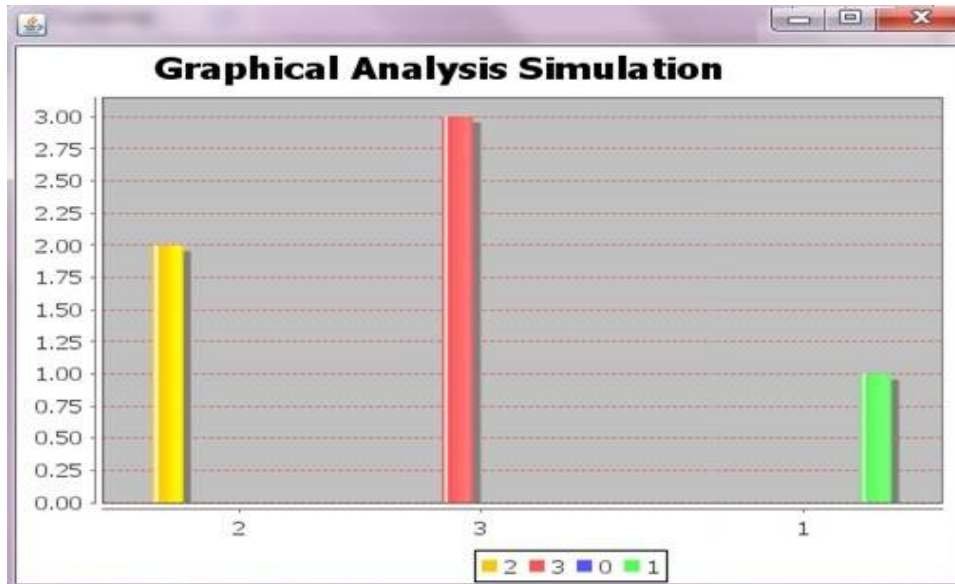


Figure 15. Bar Chart for urls in clusters

## VI. CONCLUSION

I can conclude that we have acquainted a novel system which induce the user objective oriented result. In this technique I have utilized feedback sessions to gather user search needs instead of utilizing search results or clicked URLs. Both the clicked URLs and the unclicked ones before the last click are considered as user verifiable feedbacks and considered to build input sessions. Here I have keep up the arrangement of most pertinent search results to represent need of user. I have utilized the idea of pseudo documents to outline the input sessions. This idea will make the searching simple to user. What's more, it is delivering most significant results. Experimental results on client navigate logs from a business web index show the viability of our proposed schedules. The intricacy of proposed procedure is low and I can use this framework in all fact easily. In this way by using the proposed framework client can find what he require advantageously.

## REFERENCES

- [1] R. Baeza-Yates, C. Hurtado, and M. Mendoza, "Query Recommendation Using Query Logs in search Engines," Proc. Int'l Conf. Current Trends in Database Technology (EDBT '04), pp. 588-596, 2004.
- [2] 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 (SIGKDD '00), pp. 407-416, 2000.
- [3] S. Beitzel, E. Jensen, A. Chowdhury, and O. Frieder, "Varying Approaches to Topical Web Query Classification," Proc. 30th Ann. Int'l ACM SIGIR Conf. Research and Development (SIGIR '07), pp. 783-784, 2007.
- [4] H. Cao, D. Jiang, J. Pei, Q. He, Z. Liao, E. Chen, and H. Li, "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, 2008.
- [5] H. Chen and S. Dumais, "Bringing Order to the Web: Automatically Categorizing Search Results," Proc. SIGCHI Conf. Human Factors in Computing Systems (SIGCHI '00), pp. 145-152, 2000.
- [6] C.-K Huang, L.-F Chien, and Y.-J Oyang, "Relevant Term Suggestion in Interactive Web Search Based on Contextual Data in Query Session Logs," J. Am. Soc. for Data Science and Technology, vol. 54, no. 7, pp. 638-649, 2003.
- [7] 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.
- [8] 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.
- [9] D. Kornack and P. Rakic, "Cell Proliferation without Neurogenesis in Adult Primate Neocortex," Science, vol. 294, Dec. 2001, pp. 2127-2130, doi:10.1126/science.1065467.



ISSN(Online): 2319-8753  
ISSN (Print): 2347-6710

## International Journal of Innovative Research in Science, Engineering and Technology

*(An ISO 3297: 2007 Certified Organization)*

**Vol. 5, Issue 6, June 2016**

- [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 Data Retrieval (SIGIR '05), pp. 154-161, 2005. Article in a conference proceedings.
- [11] H. Goto, Y. Hasegawa, and M. Tanaka, "Efficient Scheduling Focusing on the Duality of MPL Representatives," Proc. IEEE Symp. Computational Intelligence in Scheduling (SCIS 07), IEEE Press, Dec. 2007, pp. 57-64, doi:10.1109/SCIS.2007.357670.
- [12] Naynaneni Lavanya and E. Sandhyarani, "A Comparative Study on Privacy by Search Engines while Publishing Search Logs" International Journal of Advanced Research in Computer Science and Software Engineering, doi. 8, 2012.
- [13] Steven M. Beitzel, "Varying Approaches to Topical Web Query Classification" SIGIR 2007.
- [14] Doug Beeferman and Adam Berger, "Agglomerative clustering of a search engine query log", 2000-2010.