

# A Survey on Optimization of User Search Goals by using Feedback Session

Mr. Mahesh M. Kedari<sup>1</sup> Ms. Seefa S. Kazi<sup>2</sup> Ms. Samiksha S. Mohire<sup>3</sup> Prof. Geetanjali N. Sawant<sup>4</sup>  
<sup>1,2,3,4</sup>Department of Computer Engineering

<sup>1,2,3,4</sup>Rajendra Mane College of Engineering and Technology, Ambav, Mumbai University

**Abstract**— Data Mining refers to acquiring knowledge from large amounts of data. In the recent years the lots of surfing is done through web searching. When the information is retrieved the users clicks on particular URL, based on that click rate, ranking will be done automatically. This search engine significantly reduces the computation time required for partition of the dataset. It will also reduce the original dataset into simplified dataset. It also simplifies the data set and finds the relevant document based on users feedback. It also helps in reducing the iteration and improves the performance. Analyzing user search goal is needed to provide best result for which the user looks in the internet. Feedback sessions have been clustered to discover different user search goals for query. Pseudo-document is generated through feedback sessions for clustering. With this, the original search results are restructured. The performance of restructured search results is evaluated by classified average precision (CAP). This evaluation is used as feedback for selecting the optimal user search goals.

**Key words:** Data Mining Search Goals, Feedback Sessions, Pseudo-Documents, Classified Average Precision

## I. INTRODUCTION

The goals and objective of the information retrieval is to find the certain queries that are more relevant and the problem is to find the relevant documents. The different types of information are as follows text, audio, video, source code, application and web browser. Some static queries are posted against a static collection. The information retrieval cycle, if information retrieval cycle totally consists of five phases i.e. source collection, formulation of queries, search, selection and last result. The search process consists of the index and document collection. Indexing the documents, process the queries, evaluate similarities and find ranking and display the results are the main tasks in the information retrieval process. The documents search the queries that match the string in search box.

A Web server usually registers a log entry, for every access of a Website. It includes the URL requested IP address from which the request originated, and timestamp. Based on the Weblog records. So we have to construct the feedback session. Because log data provide information about what kind of data the users will access what kind of Web pages. This session consist of RL's and click sequence and it also focus on user search goals. Only using a feedback session we do not understand the user search goals approximately based on the feedback session construct the pseudo document for analyzing the accurate search result. This pseudo-document consist of keywords of URL's in the feedback session. This is known as enriched URL's. The enriched URL's is clustered and a pseudo document is created. Clustering is the process of collecting data into classes, so that objects within a classes have a high similarity in comparison to one another but are very dissimilar to object in other classes. After constructing the

pseudo document the Web search results are restructured based on the details collected in document. The rest of the paper contains detailed information about data mining, feedback sessions, page ranking, pseudo-document, and conclusion.

## II. PROBLEMS IN INFORMATION RENEWAL

- 1) How we can represent the documents with selected keywords?
- 2) How documents queries are correlate to calculate the weight?
- 3) Incorrect wordbook.
- 4) Ambiguous query.

## III. WEB MINING

Web mining is the use of data mining techniques to automatically discover and acquire information from Web documents and services. At Scale Unlimited we focus on the last one – acquiring value from web pages and other documents found on the web.

There are three general division of information that can be discovered by web mining:

- 1) Web activity, from server logs and Web browser activity tracking.
- 2) Web graph, from association between pages, people and other data.
- 3) Web content, for the data found on Web pages and inside of documents.

While search is the biggest web miner by far, and generates the most revenue, there are many other beneficial end uses for web mining results. A partial list includes:

- 1) Business intelligence.
- 2) Competitive intelligence.
- 3) Estimate analysis.
- 4) Events.
- 5) Product data.
- 6) Popularity.
- 7) Reputation.

## Four Steps in Content Web Mining

While acquiring the content from the internet following steps are required:

- 1) Collect: fetch or convey the content from the Web.
- 2) Parse: extract usable data from formatted data (HTML, PDF, etc).
- 3) Analyze: tokenize, weight, classify, cluster, filter, sort, etc.
- 4) Produce: turn the results of analysis into something useful (report, search index, etc).

## IV. CLASSIFICATION AND ANALYSIS, PREDICTION AND CLUSTERING

The process of grouping a set of physical or abstract objects into classes of identical or similar objects is called clustering. A cluster is a collection of data objects that are

similar to one another within the same cluster and are different to the objects in other cluster. Dissimilarities are accessed based on the attribute values describing the objects, usually distance measures are used. In this paper we use k-means clustering technique for creating pseudo documents. K-means clustering is a centroid based technique. Classification and prediction or indicator are two forms of data analysis that be used to acquire models describing important data classes or to predict future data aim. Such analysis can help provide us with a better understanding of the data at large whereas classification predicts definite labels, prediction models continuous-valued functions. It uses the pre processing technique such as data cleaning, relevance analysis, data transformation and minimization. It provide the accuracy, scalability, robustness, speed.

## V. LITERATURE REVIEW

### 1) Click Rate:

How many percent of time did the user clicks on a particular URL in a web page is given in a query.

### 2) Google Keyword Tool:

Google's keyword tool is a free online research tool to help the user for searching the appropriate keyword.

### 3) Bounce Rate:

Bounce rate is used to calculate the rate for how many times the searchers clicked on a particular webpage.

### 4) Webcap Browser Side Tool:

This tool is used to collect the relevant information from different users. It collects information based on user's interest.

### A. Page Ranking Algorithm:

The page rank algorithm anytime use in retrieve the information. Data sets available on the web can be very large determined by the rank for single page not for the whole and fill ten to hundreds of terabytes, need a large website. Based on the user click, the page rank can be form on servers. A web page contains three forms of data i.e. original page rank algorithm is regular i.e structured, unregulated i.e unstructured and semi structured data. A number of algorithms are available to make a structured data i.e.,  $PR(A) = (1-D) + D (PR(T1)/C(T1) + \dots + PR(Tn)/C(Tn))$ , one such algorithm is a fuzzy self constructing. A unstructured data can be analyzed using term frequency, document frequency, document length, text proximity. So we have to improve searching on the web by adding regulated documents. Using clustering techniques we PR (A), have to restructure the web information. We provide a PR(Ti) of pages Ti which link to page A, hierarchical classification of documents using web directories such like Google. While increasing the annual band C(Ti) on page Ti d width at ten times its average is increase three times, is a damping factor which can be set between 0's and 1's. Because of that the traffic management are important in How to increase the page rank values. The sum of the web mining pages, the rank is similar; we can add a new page to the web site. Online research tool to help the user to find the web page.

### B. Lexical Pattern Extraction Algorithm:

In information improvement, one of the main problems is to retrieve a set of documents that is syntactically related to a given user query. We need to identify the numerous semantic we cannot choose ranking functions randomly. In order relations that lie between two given words. Efficient to ascertain the correct ranking function, we use the retrieval information is based on words. The information retrieval is concept called query similarity. It has two distinct problem when the words are same. On the other approaches, query condition similarity and query result hand, the snippets /fragment returned by a search engine for the similarity.

### C. Virtual Feature Updating Algorithm:

VFs capture the words. Fragment/snippets contain a window of text selected from high level semantics imposed by multiple users. Document that includes the queried words, Fragments are semantics concepts are learned information retrieval use for user while searching, users can learn the action in two ways: In breadth and in depth. The fragments are divided without opening the URLs. Using fragment virtual feature technique uses the both types of learning, is productive it is helpful to download the useful document from to produce the desired result in real-world applications. In the web it is time consuming task when the document is large. Breadth is learned from one session in a single query.

### D. Construction of Virtual Documents:

The virtual document, for in-breadth learning, the user can utilize the document terms are acquired from its detailed graph. System to increase the weight of the conception in the virtual document consists of terms single query session. For in-depth learning, is based extracted from its description graph: For each entity on previous or past information retrieval with multiple users, the identified by a URL in this graph, extract its local name and system automatically utilize the long-term knowledge and label; for each literal in this graph, acquire its phrasal form. it is used to search the relevant images in the future

### E. Index Based Threshold Algorithm:

The algorithm is "Customer," "type," "Restriction," "on Property," based on two accesses: Sorted and random access. Sorted "studies at," "all Values From," and "University." It is access retrieves the tuples information based on best to property names are also included so that the decreasing order of their attribute value. It maintains a system can support more opposite keyword queries, e.g. buffer and allows using a stopping condition due to the swrc: Student can be retrieved by the keyword doubt detection of final Top K tuples before processing the all "subclass of person." tuples.

### F. User and Doubt Similarity Model:

When the multiple the query is no more than the tuples in the Top-K bufferusers in the organization, the two users querying the equal with the lowest similarity, the algorithm successfully set of queries. The ranking functions have been derived terminates on browser choices. The user similarity between the two users can be expressed as the average similarity between queries used by both users.

### G. Filtering Algorithm:

The algorithm stocks the data in the ranking functions. The goal of the similarity is to complex points in a Kd - Tree. Kd-tree is determined the ranking functions derived from the similar binary tree having nodes and leaf. Each node of a kd -Tree doubt from the similar user is called cell, which contains closed box. The more than when the user in the system enters the query in the one point in a cell is called a

leaf. The points in the cell are browser, the results are obtained. Let be the ranking partitioned into one side. The remaining cells are children functions derived from the each individual query of the original cell. There are a number of ways which are already present. The same query having different ranking functions.

## VI. COMPARATIVE STUDY

Study Sr. no	Titles	Technique/Methods	Advantage	Disadvantage
1.	Automatic Credential of User Goals in Web Search	User click behavior and Anchor link distribution	Using goal identification task to achieve approximately 90% of accurate results	Probably-based dataset
2.	Query-Sets: Using Implied Feedback and Query Patterns to Organize Web Documents	Non supervised tasks	Improve the quality approximately as 90%	A broader comparison with online directory
3.	Learn from Web Search Logs to Construct Search Results	Commercial search engine log data and clustering	Better result construction and meaningful labels	descriptive feedback information from user
4.	Generating Query Substitutions	Query combination algorithm	Increase coverage and effectiveness	Machine translation techniques
5.	Learning Query Decided from Consistent Click Graphs	Semi-supervised click graph	Improve classification performance	Impact of concept queries and divide query classification
6.	Varying Approaches to Popular Web Query Classification	Pre vs. post retrieval classification	QC is outperforms link a document anatomy as 48%	multiple approaches to improve performance
7.	Text-Aware QS by Mining Click-Through and Session/Term Data	Offline model learning and online QS step, concept sequence is additional	Analysis and quality of suggestions	Larger analysis area

## VII. CONCLUSION

As per the references mentioned techniques above techniques deals with information retrieval process on internet.

Web technology and its usage continue to grow, so grows the opportunity to analyze Web data and acquire all manner of useful knowledge from it. The previous few years have seen the emergence of Web mining as a rapidly growing area, due to the efforts of the research association as well as various institute that are practicing it. In this paper, we propose a innovative access for user search goals using feedback session and pseudo document. First we construct a feedback session to analysis the user search goal from the Weblog record. It cannot provide the approximate result. So that we can introduce the pseudo documents to provides the approximate results. By using pseudo documents we have to restructure the Web search results.

## ACKNOWLEDGEMENT

Project is never complete without the guidance of those expert who have already traded this past before and hence become skillful of it and as a result, our leader. So we would like to take this opportunity to take those entire individual they have helped us in vocalizing this project. We express our deep gratitude to our project guide Prof. G.N. Sawant for providing timely assessments' to our query and guidance that she gave owing to her experience in this field for past many. We would also take this opportunity to thank our

project co-ordinate Prof. G. N. Sawant for her guidance on selecting this project and also for providing us all this details in proper presentation of this project. We extend our sincere appreciation to our entire Professor from RAJENDRA MANE COLLEGE OF ENGINEERING & TECHNOLOGY for their valuable guidance and important tips during the designing of the project. Their contributions have been valuable in so many ways, we find it difficult to acknowledge of them individual. We also grateful to our HOD Prof. L. S. Naik, for providing help directly and indirectly through various difficulties in our project work.

## REFERENCES

- [1] Beitzel. S, 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.
- [2] Baeza-Yates. R, 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.
- [3] [3] Beeferman. D 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.
- [4] Cao. H, 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] Chen. H 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] Huang C.K, L.-F Chien, and Y.-J Oyang, "Relevant Term Suggestion in Interactive Web Search Based on Contextual Information in Query Session Logs," J. Am. Soc. for Information Science and Technology, vol. 54, no. 7, pp. 638-649, 2003.
- [7] Jones. R, B. Rey, O. Madani, and W. Greiner, "Generating Query Substitutions," Proc. 15th Int'l Conf. World Wide Web (WWW '06), pp. 387-396, 2006.
- [8] Joachims. T, 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.
- [9] Joachims. T, "Evaluating Retrieval Performance Using Clickthrough Data," Text Mining, J. Franke, G. Nakhaeizadeh, and I. Renz, eds., pp. 79-96, Physica/Springer Verlag, 2003.
- [10] Joachims. T, "Optimizing Search Engines Using Clickthrough Data," Proc. Eighth ACM SIGKDD Int'l Conf. Knowledge Discovery and Data Mining (SIGKDD '02), pp. 133-142, 2002.

