

# Information Retrieval Using Keyword Query Analysis

A.V.Akhare<sup>1</sup>, Dr. M. A. Pund<sup>2</sup>

<sup>1</sup>Dept of CSE

<sup>2</sup>Professor, Dept of CSE

<sup>1,2</sup>PRMIT & R Badnera, Amravati, India

**Abstract-** Information Retrieval is concerned with indexing and retrieving documents including information relevant to a user's information need. Relevance Feedback (RF) is a class of effective algorithms for improving Information Retrieval (IR) and it consists of gathering further data representing the user's information need and automatically creating a new query. In this paper, we propose a class of RF algorithms inspired by quantum detection to re-weight the query terms and to re-rank the document retrieved by an IR system. These algorithms project the query vector on a subspace spanned by the eigenvector which maximizes the distance between the distribution of quantum probability of relevance and the distribution of quantum probability of non-relevance. The experiments showed that the RF algorithms inspired by quantum detection can outperform the state-of-the-art algorithms.

**Keywords-** Information retrieval, quantum mechanics, relevance feedback, quantum detection

## I. INTRODUCTION

Finding relevant document is one of the hardest tasks. Information retrieval (IR) is the activity of obtaining information resources relevant to an information need from a collection of information resources. Searches can be based on full-text or other content-based indexing. Automated information retrieval systems are used to reduce what has been called "information overload". Many universities and public libraries use IR systems to provide access to books, journals and other documents. Web search engines are the most visible IR applications. An information retrieval process begins when a user enters a query into the system. Queries are formal statements of information needs, for example search strings in web search engines. In information retrieval a query does not uniquely identify a single object in the collection. Instead, several objects may match the query, perhaps with different degrees of relevancy. An object is an entity that is represented by information in a content collection or database. User queries are matched against the database information. However, as opposed to classical SQL queries of a database, in information retrieval the results returned may or may not match the query, so results are typically ranked. This ranking of results is a key difference

of information retrieval searching compared to database searching. Depending on the application the data objects may be, for example, text documents, images, audio, mind maps or videos. Often the documents themselves are not kept or stored directly in the IR system, but are instead represented in the system by document surrogates or metadata. Most IR systems compute a numeric score on how well each object in the database matches the query, and rank the objects according to this value. The top ranking objects are then shown to the user. The process may then be iterated if the user wishes to refine the query. We propose to investigate the problem of keyword query routing for keyword search over a large number of structured and Linked Data sources. Routing keywords only to relevant sources can reduce the high cost of searching for structured results that span multiple sources. Information retrieval (IR) has experienced huge growth in the past decade as increasing numbers and types of information systems are being developed for end-users. The incorporation of users into IR system evaluation and the study of users information search behaviours and interactions have been identified as important concerns for IR researchers [11]. The proposition that IR systems are fundamentally interactive and should be evaluated from the perspective of users is not new. An IR system addresses the problems caused by query ambiguity by gathering additional evidence that can be used to automatically modify the query. Usually a query is expanded because the queries are short and it cannot exhaustively describe every aspect of the user's information need; however, some irrelevant documents may be retrieved or relevant documents may also be missed when a query is not short. [10] The automatic procedure that modify the user's queries is known as Relevance Keyword query routing (RF); some relevance assessments about the retrieved documents are collected and the query is expanded by the terms found in the relevant documents, reduced by the terms found in the irrelevant documents or reweighted using relevant or irrelevant documents.

## RELATED WORK

Author : Ingo Frommholz Paper title: Supporting Poly representation in a Quantum-inspired Geometrical Retrieval Framework Proposed Methodology: Geometrical retrieval framework inspired by quantum mechanics can be extended to

support poly representation . This system was unable to show that the well-motivated algorithms perform significantly better than the simple algorithms. Luis M. de Campos Implementing Relevance Feedback in the Bayesian Network Retrieval Mode Publication: JOURNAL OF THE AMERICAN SOCIETY FOR INFORMATION SCIENCE AND TECHNOLOGY, 54(4):302–313,2003.present an approach for relevance feedback in the Bayesian Network Retrieval (BNR) model. It is very difficult to compare feedback methods when the retrieval engines are differentGeneral purpose propagation algorithms can't be applied due to efficiency considerations. Automatically learning the relationships among terms could imply that some relationships are not strong enough. Retrieval effectiveness could be damaged. If the number of terms is very high, the learning stage could be time consuming.M. Shanmugham BEARINGSTIMULATEDALGORITHMS INSPIRED BY QUANTUM DETECTION Publication: International Journal of Current Trends in Engineering & Research (IJCTER) e-ISSN 2455–1392 Volume 2 Issue 7, July 2016 pp. 249 – 255 Proposed Methodology: Present a class of RF algorithms inspired by the quantum detection has been proposed to re-weight query terms by projecting the query vector on the subspace represented by the eigenvector . Explicit RF , Pseudo RF and Implicit RF is based on observations that are proxies of relevance. The main problem with proxies is that they are not necessarily reliable indicators of relevance and thus should be consider noisy. These systems do no rely on only non-retrieval technology. CLAUDIO Automatic Query Expansion in Information Retrieval. ACM Computing Surveys, Vol. 44, No. 1, Article 1, Publication date: January 2012.

**Motivation**Information Retrieval is the process of obtaining relevant information from a collection of informational resources. It does not return information that is restricted to a single object collection but matches several objects which vary in the degree of relevancy to the query. So, we have to think about what concepts IR systems use to model this data so that they can return all the documents that are relevant to the query term and ranked based on certain importance measures. These concepts include dimensionality reduction, data modeling, ranking measures, clustering etc. The tools that IR systems provide would help you get your results faster. So, while computing the results and their relevance, programmers use these concepts to design their system, think of what data structures and procedures are to be used which would increase speed of the searches and better handling of data

**.Aim**To maintain the collection of documents according to different user search, To find query-document or document-document similarity. The reduction is not really substantial, To measure the performance relevance judgments more

accurately and more quickly. Users can identify more relevant documents for each query, while at the same time make fewer mistakes, to Find the document according to content of the documents, to Implement the concept of relevant document suggestion.

## II. LITERATURE SURVEY

**InformationRetrival** Information Retrieval is the process of obtaining relevant information from a collection of informational resources. Information retrieval (IR) is the activity of obtaining information resources relevant to an information need from a collection of information resources. Searches can be based on full-text or other content-based indexing. Information retrieval is the science of searching for information in a document, searching for documents themselves, and also searching for metadata that describe data, and for databases of texts and other types of document.

**Data Mining**Data mining is the process of discovering patterns in large data sets involving methods at the intersection of machine learning, statistics, and database systems.[1] It is an essential process where intelligent methods are applied to extract data patterns. It is an interdisciplinary subfield of computer science.[4] The overall goal of the data mining process is to extract information from a data set and transform it into an understandable structure for further use. Aside from the raw analysis step, it involves database and data management aspects, data pre-processing, model and inference considerations, interestingness metrics, complexity considerations, post-processing of discovered structures, visualization, and online updating. Data mining is the analysis step of the "knowledge discovery in databases" process, or KDD.

**Keyword Query**Keyboard Query is the keyword entered by user for searching purpose in the document. using keyword query searching document is relevant to user. The entered keyword is matched with the content of document and relevant document is retrieved

## III. PROPOSED SYSTEM ANALYSIS & DESIGN

**Analysis** System analysis is a problem solving technique that decomposes a system into its component pieces for the purpose of the studying how well those component parts work and interact to accomplish their purpose. From the above literature survey the problem is identified and formulated as follows.

**Problem Definition**An IR system addresses the problems caused by query ambiguity by gathering additional evidence

that can be used to automatically modify the query. Besides negativeness and positiveness, the RF algorithms can be classified according to the way the relevance assessments are collected. Keyword query routing may be explicit when the user explicitly tells the system what the relevant documents and the irrelevant documents are top-ranked documents are considered as relevant documents, or it is implicit when the system monitors the user's behavior and decides what the relevant documents and the irrelevant documents are according to the user's actions.

**Proposed System**We are going to propose a system using which the user can easily get the relevance document. When the user enter the query for search the document, then it directly compare within the data of the document file. So the relevant document will found by the system. We are also working to add feature, the system will recommend the keyword to the user for getting the best result or document.Objectives of this work is To maintain the collection of documents according to different user search. To find query-document or document-document similarity. The reduction is not really substantial, To measure the performance relevance judgments more accurately and more quickly, Users can identify more relevant documents for each query; while at the same time make fewer mistakes, also To Find the document according to content of the documents.

**Algorithm Rocchio's Algorithm**

The Rocchio algorithm is based on a method of relevance feedback found in information retrieval systems which stemmed from the SMART Information Retrieval System. Like many other retrieval systems, the Rocchio feedback approach was developed using the Vector Space Model. The algorithm is based on the assumption that most users have a general conception of which documents should be denoted as relevant or non-relevant. Therefore, the user's search query is revised to include an arbitrary percentage of relevant and non-relevant documents as a means of increasing the search engine's recall, and possibly the precision as well. Explicit Relevance FeedbackIt also called as Term relevance feedback. The system will suggest the term which types of term the user should add in search. Implicit Relevance FeedbackIt will find out the frequently search document easily.Jun Miao [26] we study how to incorporate proximity in- formation into the Rocchio's model, and propose a proximity- based Rocchio's model, called PRoc, with three variants. In our PRoc models, a new concept (proximity-based term frequency, ptf) is introduced to model the proximity information in the pseudo relevant documents, which is then used in three kinds of proximity measures. Experimental results on TREC collections show that our proposed PRoc

models are effective and generally superior to the state-of-the-art relevance keyword query routing models with optimal parameters. A direct comparison with positional relevance model (PRM) on the GOV2 collection also indicates our proposed model is at least competitive to the most recent progress.

**Proposed System Architecture**We are going to propose a system using which the user can easily get the relevance document. When the user enter the query for search the document, then it directly compare within the data of the document file. So the relevant document will found by the system. We are also working to add feature, the system will recommend the keyword to the user for getting the best document.

**Data flow diagram**This is stage of the project when the theoretical design is turned out in working system. Thus it can be considered to be the most critical stage in achieving a successful new system and in giving the user, confidence that the new system will work and be effective. The framework involves careful planning, investigation of existing system and its constraints on implementation, designing of methods to achieve goal. In the first step users have to register with information like name, password and type of the user. Then with the help of userid and password user can login. Enter query into IR system and search document IR system shows relevant document that user want. Suppose user enter query in such a way that IR system has no document related to query content in that situation IR system match the query content with the content inside the document and return relevant document that user want.

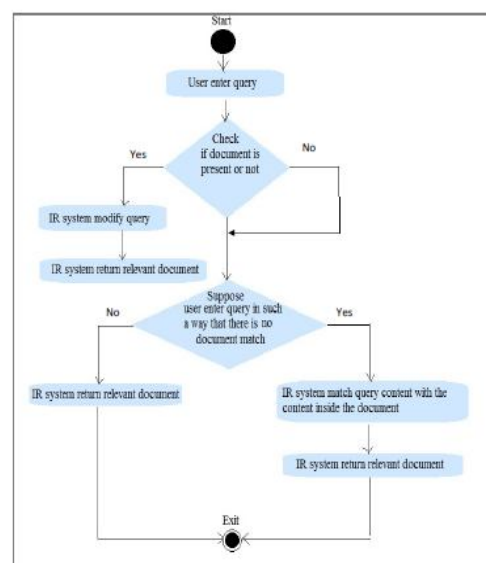


Figure 1: Data-Flow Diagram

**Sequence Diagram** A sequence diagram is a type of interaction diagram. In particular it shows the objects participating in the interaction by their lifelines and messages that they are exchanged.

- Step 1:** Create an account: User creates an account whose user id and password is stored in database.
- Step 2:** Log in: User will login through username and password in system. System will check user id and password from the database.
- Step 3:** Gives authentication: The authorized user will be authenticated.
- Step 4:** User need some information so user enter simple query into IR system.
- Step 5:** IR system display Relevant document that user need.
- Step 6:** Suppose user enter long query.
- Step 7:** IR system modify query and display relevant document.
- Step 8:** Suppose user enter query in such a way that there is no document match.
- Step 9:** IR system match query content with content inside the document and return relevant document.
- Step 10:** Log out: User will log out.

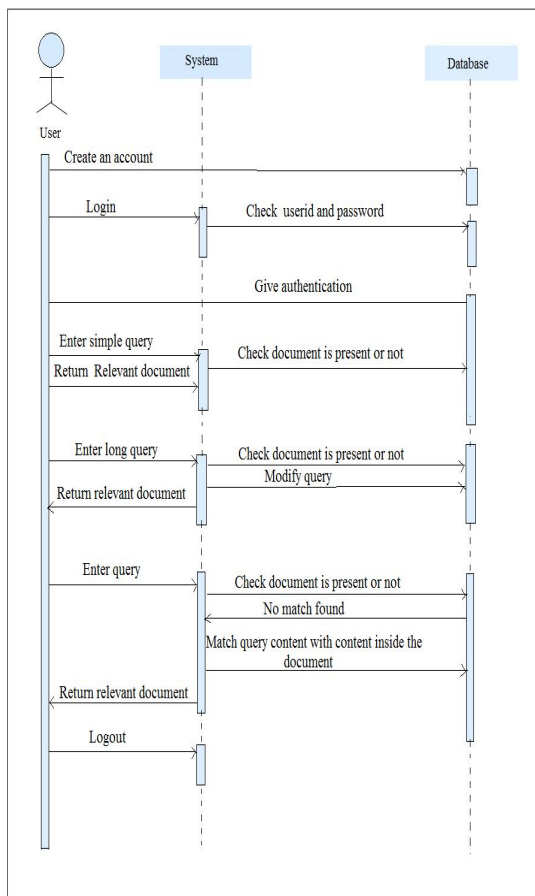


Fig no 2. Sequence Diagram

#### IV. CONCLUSION

Main objective of relevant information retrieval using keyword query is to To measure the performance relevance judgements more accurately and more quickly also To maintain the collection of documents according to different categories and according to different user search. The result given by system is according to the priority or the sequence of the more no of word occurred in each document , the document having more no of word entered during keyword search is placed at the first place of search result along with the time required for each document. Relevance keyword query routing can go through one or more iterations of this sort. The process exploits the idea that it may be difficult to formulate a good query when you don't know the collection well, but it is easy to judge particular documents, and so it makes sense to engage in iterative query refinement of this sort. In such a scenario, relevance keyword query routing can also be effective in tracking a user's evolving information need: seeing some documents may lead users to refine their understanding of the information they are seeking. The user submits a query into IR system. IR system return both relevant and irrelevant documents so the automatic procedure that modify the user's queries is known as query routing; IR system return relevant documents that user need with processing time require for each document. Some relevance assessments about the retrieved documents are collected and the query is expanded by the terms found in the relevant documents, reduced by the terms found in the irrelevant documents or reweighted using relevant or irrelevant documents.

#### REFERANCES

- [1] Luis M. de Campos, Juan M. FernándeZ-Luna ,Juan F. Huete "Implementing Relevance Keyword query routing in the Bayesian Network Retrieval Model ",*JOURNAL OF THE AMERICAN SOCIETY FOR INFORMATION SCIENCE AND TECHNOLOGY*, 54(4):302–313, 2003.
- [2] X. Tian, L. Yang, J. Wang, X. Wu, and X. Hua, "Bayesian visual reranking," *IEEE Trans. Multimedia*, vol. 13, no. 4, pp. 639–652, 2011
- [3] Diane Kelly, "Methods for Evaluating Interactive Information Retrieval Systems with Users ", *Foundations and Trends in Information Retrieval* Vol. 3, Nos. 1–2 (2009) 1–224 c 2009.
- [4] Shuqin Liu, Jinye Peng, "A Novel Image Retrieval Algorithm Based on Adaptive Weight Adjustment and Relevance Keyword query routing ",*JOURNAL OF COMPUTERS*, VOL. 9, NO. 11, NOVEMBER 2014.
- [5] Ingo Frommholz, Birger Larsen, Benjamin Piwowarski, Mounia Lalmas, Peter Ingwersen, Keith van Rijsbergen,

- "Supporting Polyrepresentation in a Quantum-inspired Geometrical Retrieval Framework". IiX 2010, August 18–21, 2010.
- [6] K. Collins-Thompson, P. N. Bennett, R. W. White, S. de la Chica, and D. Sontag, "Personalizing web search results by reading level", in Proc. 20th ACM Int. Conf. Inf.Knowl.Manage., 2011, pp. 403–412.
- [7] W. Liu, G. Hua, and J. Smith, "Unsupervised One-Class Learning for Automatic Outlier Removal", in Proc. IEEE Int. Conf. Comput. Vis. Pattern Recognit., 2014, pp. 3826–3833
- [8] AblimitAji, Yu Wang, EugeneAgichtein, EvgeniyGabrilovich, "Using the Past To Score the Present: Extending Term Weighting Models Through Revision History Analysis". CIKM'10, October 26–30, 2010.
- [9] Yuanhua, Cheng Xiang Zhai, Wan Chen. "A Boosting Approach to Improving Pseudo-Relevance Keyword query routing." SIGIR'11, July 24–28, 2011.
- [10] G. Salton and C. Buckley. *Improving retrieval performance by relevance keyword query routing*. Journal of the American Society for Information Science. 41. 4. pp 288-297. 1990.