

High-Dimensional Information Filtering Using Query Reorganization Algorithms

¹Mr.J.Srinivasan, ²K.Thavamani

¹Assistant Professor, ²M.Phil (CS) Research Scholar

^{1,2}Department of Computer Science and Applications,

^{1,2}Adhiparasakthi College of Arts and Science (Autonomous), G.B.Nagar, Kalavai -632506, Vellore (District)

Corresponding Author: Mr.J.Srinivasan

Abstract : The information filtering system used to publish information needs of clients and relevant information to a server with repeated queries. Indexing schemes holding quick meet of incoming information with query database are employed to the servers to publish information in an efficient way. The indexing schemes include (i) main-memory trie-based data structures that cluster similar queries by capturing common elements between the queries, and (ii) efficient filtering mechanisms that utilize query clustering to ensure high throughput and low filtering times. But indexing schemes are sensitive to the query insertion order and cannot adapt to an evolving query workload, and create filtering work over time. To overcome the drawbacks of the indexing schemes, we present an adaptive trie based algorithm which gives better current techniques by relying on query statistics to rearrange the query database. The algorithm provides efficient filtering performance as determining factor by showing nature of the constructed tries, rather than compactness of the constructed tries. The algorithm does not depend on the order of insertion of queries in the database, and manages to cluster queries even when clustering possibilities are limited, so that the algorithm ensures improved filtering time than other systems. We will also demonstrate that the algorithm can be easily extensible to multi-core machines.

Keywords -Information filtering, user profiles and alert services, indexing methods, dissemination.

Date of Submission: 21-07-2018

Date of acceptance: 6-08-2018

I. INTRODUCTION

Information filtering (IF) applications (known as information dissemination or publish/subscribe), such as news alerts, weather monitoring, and stock quotes, have gained popularity in recent times. These applications assist users to cope with the information avalanche and the cognitive overload associated with the information avalanche. Data of interest is mostly textual in some applications like news alerts, digital libraries, or RSS feeds where users of the applications express their needs to a server using information retrieval languages (e.g. Boolean combinations of keywords or text excerpts under the Vector Space Model – VSM and submit continuous queries (or profiles) so as to subscribe to newly appearing documents that will satisfy the query conditions.

Then the server will be responsible for notifying the subscribed users automatically whenever a new document matching user's information needs is published, where publishers can be news feeds, digital libraries, or even users who post new items to blogs, social media, and Internet communities. This functions of IF is different from information retrieval (IR) applications like search engines. When a query is posed in IR, a single search is executed and the current matching data items are presented to the user but in IF the server indexes the user queries rather than the data and evaluates newly published data items for the stored continuous queries. In detail, information filtering having following problems: a database (DB) of continuous queries that reside on a server and an incoming document (d) are set to retrieve all queries q2DB that match. These filtering problem needs to be solved efficiently, since servers are expected to handle millions of user queries and high rates of published documents. Efficiency issues were identified by many researchers that proposed tree and trie-based algorithms for supporting fast filtering under various data models (e.g., flat attribute-based, semi-structured XML) and query languages (e.g., Boolean, VSM), both for main-memory and secondary storage. But all these approaches use a greedy clustering method which is sensitive to the insertion order of submitted queries and do not consider that an evolving query workload might require the reorganization of the query database to ensure efficient filtering performance.

II. LITERATURE SURVEY

2.1 Introduction to Information Retrieval

In this paper, few new indexing methodologies and applications in Information Retrieval (IR) has been presented. We also introduced some new algorithms with high coverage of IR applications. Main strategy is introducing and evaluating Information Retrieval basic applications and modulation. Some future directions in IR methodologies and evaluations are also included as other subjects and focuses on this paper.

2.2 Batched Processing for Information Filters

This paper describes batching, a novel technique in order to improve the throughput of an information filter (e.g. message broker or publish & subscribe system). Rather than processing each message individually, incoming messages are reordered, grouped and a whole group of similar messages is processed. This paper presents alternative strategies to do batching. Extensive performance experiments are conducted on those strategies so as to compare their tradeoffs.

2.3 Index Structures for Information Filtering under the Vector Space Model

A retrieval model is effectively established by author following a study of what data structures and algorithms can be used to efficiently perform large-scale information filtering under the vector space model. The idea of the standard inverted index to index user profiles is applied in the retrieval model, where they select only the significant ones to index instead of indexing every term in a profile by setting up an alternative to the standard inverted index. They evaluate their performance and show that the indexing methods require orders of magnitude fewer I/Os to process a document than when no index is used. They also show that the proposed alternative performs better in terms of I/O and CPU processing time in most cases.

2.4 Document Filtering With Inference Networks

We develop a new approach for text document filtering based on automatic construction of filtering profiles using Bayesian inference network learning. Based on probability theory, Bayesian inference networks offer a suitable framework to harness the uncertainty found in the nature of the filtering problem. To learn the networks effectively, we introduce three different techniques for discretization, where these techniques automatically obtain better features of high predictive power from the training document content. Our approach does not need to know in advance the subject or content of documents as well as the information needs expressed as topics. A series of experiments on a set of topics were conducted on two large-scale realworld document corpora. The empirical results demonstrate that our Bayesian inference network learning with advanced discretization ensures better performance over the simple naive Bayesian approach.

III. SYSTEM STUDY

3.1 Existing System:

- Efficiency problems were known by several researchers that proposed tree and trie-based algorithms for supporting quick filtering below numerous information models (e.g., flat attribute-based, semi-structured XML) and query languages (e.g., Boolean, VSM), each for main-memory and external storage.
- However, all these approaches use a greedy cluster technique that is sensitive to the insertion order of submitted queries and will not take into account that an evolving query work would possibly need the reorganization of the query information to realize efficient filtering performance.

3.2 Disadvantage Existing System:

- The problem of data filtering also outlined as follows: Information of continuous queries that reside on a server and an incoming document are given to retrieve all queries that match document.
- Filtering efficiency is needed to be achieved, since servers are expected to handle numerous user queries and high rates of revealed documents.
- Limitations of the proposed family of algorithms include (i) reduced efficiency on limited query vocabularies and/or very short continuous queries, (ii) increased memory usage for indexing queries with disjunctions as the different disjunctions need to be split and indexed at different tries, and (iii) corpus-dependent parameter/algorithm setup.

3.3 Proposed System:

- The main aim behind the proposed method is to use tries to capture common components of queries, equally but, the key variations with these approaches lie (i) the gathering and utilization of statistics on the importance of keywords within the indexed queries, (ii) the reorganization of the query information according each to word and query importance, and (iii) the demonstration that the character of the trie forest is additional necessary than its compactness once it involves filtering efficiency.
- Apparently, all previous works were aiming at minimizing the dimensions of the trie forest, since there was an implicit conjecture that a little forest would lead to lower filtering times as a result of less node visits.

3.4 Advantage Of Proposed System:

- ✓ Here we use linguistic method of concepts since linguists read through large amounts of data, including texts, audios, and videos, they are trained to search for essential information among piles of data.
- ✓ Through this process, linguists gain intuition as to where and how to approach information. To support continuous queries that are comprised of conjunctions of keywords linguistic method used and it may be used as a basis for query languages that support not only basic Boolean operators, but also more complex constructs, such as proximity operators and attributes.
- ✓ Here we use an efficient Boolean filtering service (like Vector Space Model Queries) as a valuable addition to any text filtering setup. It is mostly used for index terms.

IV. DETAILED DESIGN

4.1 Implementation

Our information filtering method uses linguistically motivated concepts, such as words, to support continuous queries that are comprised of conjunctions of keywords and may be used as a basis for query languages that support not only basic Boolean operators, but also more complex constructs, such as proximity operators and attributes. We believe that offering an efficient Boolean filtering service (possibly alongside a more popular model like VSM) is a valuable addition to any text filtering setup. Boolean IR/IF is still the model of choice of advanced users that want total control of their results and is widely supported in systems of major stakeholders like Google's advanced search/alert mechanisms. Such systems, that are meant to cope with a high workload and are designed for efficiency, are possible applications for our work.

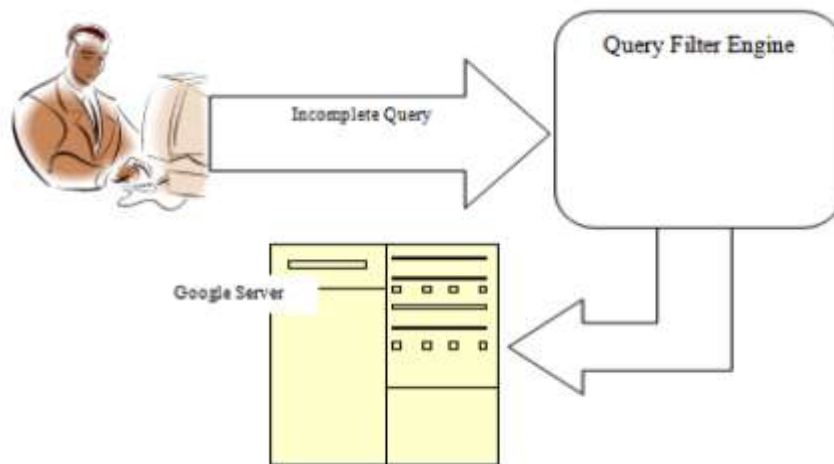


Figure 1: Information Filtering Method

4.2 User profiles and alert services

The first algorithm to identify the importance of query insertion order and its influence in the filtering time was Algorithm RETRIE. Algorithm RETRIE introduces the concept of query relocation, identifies poorly indexed queries and re-indexes them in better positions, and achieves a limited form of re-organization in the query database.

Query insertion order Influences initial creation of the tries. Contrary to the aforementioned approaches, our proposal is the first in the literature that emphasizes on the reorganization of the query database and addresses the issue of query insertion order.

4.3 Indexing methods

Other approaches included statistical filtering systems, such as Latent Semantic Indexing used to filter incoming documents and that utilizes network-based profile representations to better identify user interests and cope with the curse of dimensionality in VSM. Adaptive filtering focuses also on profile effectiveness and considers the adaptation of VSM queries and their dissemination thresholds. In order to enhance user information discovery, developed a novel statistical latent class model that applies user/item grouping to deliver better content recommendations/predictions. Moreover, sophisticated user profiling has also been used to promote personalized IR systems that focus on improving retrieval effectiveness.

4.4 Dissemination Algorithms

RETRIE and TREE algorithms ensure low sensitivity to query database size, query length, and document size while presenting. Although Algorithm STAR-HR is designed for query databases that are unfocused and cover thematically a wide variety of topics, which also performs quick filtering for both focused query databases with restricted vocabularies and real-life query logs.

Our experiments showed that Algorithm STAR-HR outperforms its competitors in terms of filtering time for various document sizes. Insertion and re-organization times for STAR-HR are also efficient and quicker than its competitors due to the placement of rare words near trie roots.

V. Result

5.1 Screenshot



Figure 2: Search Activity by Info Filtering



Figure 3: Search result by info filtering



Figure 4: View All Activities

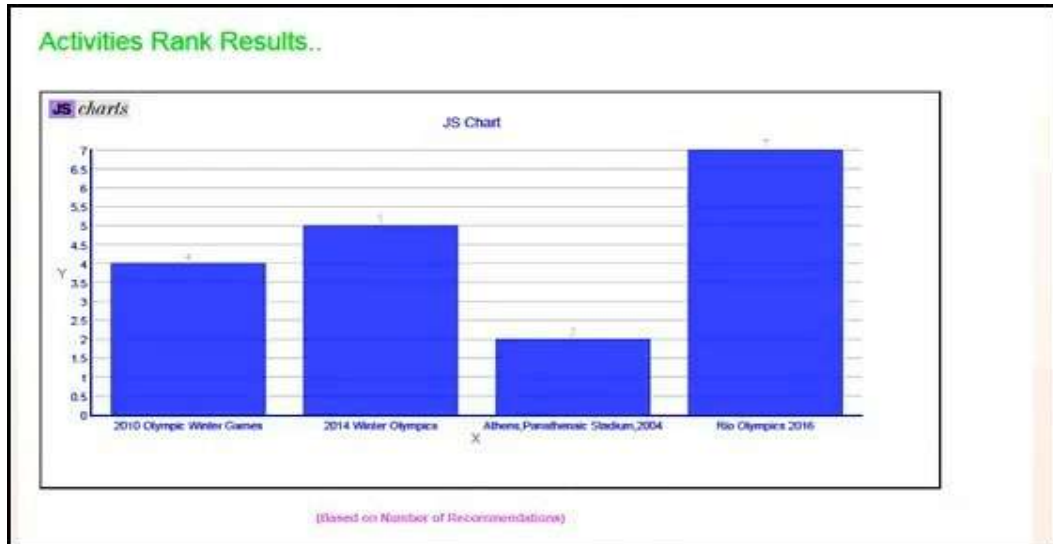


Figure 5: Activity Rank Results

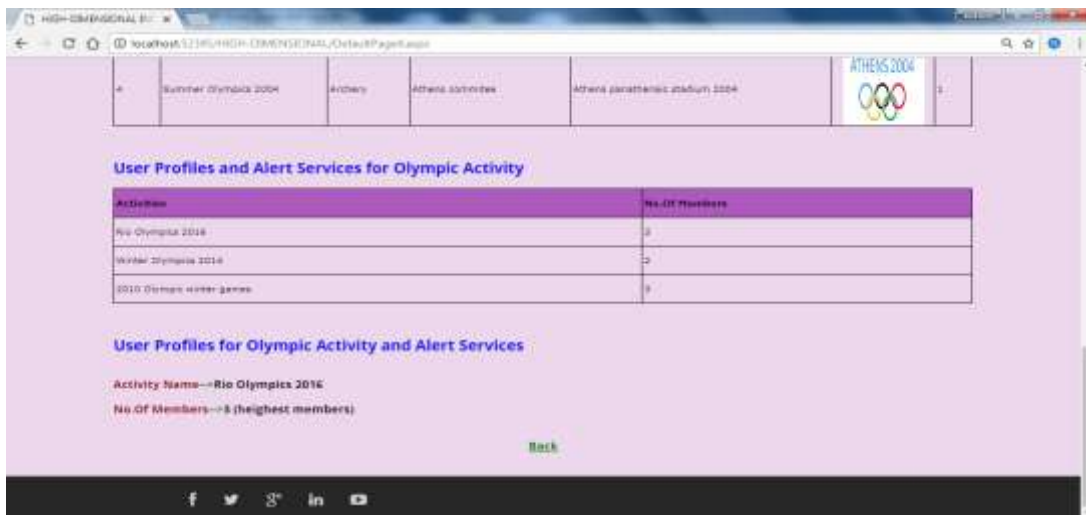


Figure 6: User Activities Details



Figure 7: View and Edit Delete Activities



Figure 8: View user Friends Activities

VI. Conclusion & Future Enhancement

6.1 Conclusion

Here we create a novel algorithm for implementing these indexing schemes which supports Vector Space Model queries. However, this data structure is designed for arithmetic and string operations and is not applicable in textual IF. Here we experimentally evaluate different rearrangement strategies and showcase their effect in filtering efficiency using two different real-world datasets and both synthetic and real query sets.

6.2 Future Enhancement

Alternatively, the current reorganization algorithm bases its decisions purely on the user accesses. It would be very interesting to study the possibility of incorporating knowledge of the Web pages by Web content mining to increase the effectiveness and reliability of the reorganization algorithm.

REFERENCES

- [1]. C. Manning, P. Raghavan, and H. Schutze, Introduction to Information Retrieval. Cambridge University Press, 2008.
- [2]. P. Fischer and D. Kossman, —Batched Processing for Information Filters, I ICDE, 2005.
- [3]. C. Tryfonopoulos, M. Koubarakis, and Y. Drougas, —Filtering Algorithms for Information Retrieval Models with Named Attributes and Proximity Operators, I ACM SIGIR, 2004.
- [4]. J.Savithri, H.Inbarani, —Comparative Analysis Of K-Means, PSO-K-Means, And Hybrid PSO Genetic K-Means For Gene Expression Data, International Journal of Innovations in Scientific and user database server register login Information filtering user profiles and alert services indexing methods respons to user International Journal of Pure and Applied Mathematics Special Issue 307 Engineering Research (IJISER), Vol.1, no.1, pp.43-50, 2014.
- [5]. Information filtering and query indexing for an information retrieval model, I ACM TOIS, 2009.
- [6]. T. Yan and H. Garcia-Molina, Index structures for selective dissemination of information under the boolean model, I ACM TODS, 1994.
- [7]. J. Yochum, A High-Speed Text Scanning Algorithm Utilising Least Frequent Trigraphs, I IEEE SNDC, 1985.
- [8]. T. Bell and A. Moffat, The Design of a High Performance Information Filtering System, I ACM SIGIR, 1996.
- [9]. T. Yan and H. Garcia-Molina, Index Structures for Information Filtering under the Vector Space Model, I ICDE, 1994.
- [10]. J. Callan, Document Filtering With Inference Networks, I ACM SIGIR, 1996.
- [11]. W. Rao, L. Chen, S. Chen, and S. Tarkoma, Evaluating continuous top-k queries over document streams, I World Wide Web, 2014.
- [12]. M. Franklin and S. Zdonik, Data in Your Face: Push Technology in Perspective, I SIGMOD Record, 1998.
- [13]. M. Altinel, D. Aksoy, T. Baby, M. Franklin, W. Shapiro, and S. Zdonik, DBIS-toolkit: Adaptable Middleware for Large-scale Data Delivery, I in ACM SIGMOD, 1999.
- [14]. F. Fabret, H. A. Jacobsen, F. Lirbat, J. Pereira, K. A. Ross, and D. Shasha, Filtering algorithms and implementation for very fast publish/subscribe systems, I ACM SIGMOD, 2001.
- [15]. B. Nguyen, S. Abiteboul, G.Cobena, and M. Preda, Monitoring XML Data on the Web, ACM SIGMOD, 2001.
- [16]. A. Campailla, S. Chaki, E. Clarke, S. Jha, and H. Veith, Efficient Filtering in Publish Subscribe Systems Using Binary Decision Diagrams, I in ICSE, 2001.

Mr.J.Srinivasan "High-Dimensional Information Filtering Using Query Reorganization Algorithms "International Journal of Engineering Science Invention (IJESI), vol. 07, no. 8, 2018, pp 68-73