

Implementing High Performance Retrieval Process by Max-Score Ranking

U.Vignesh¹, M.Sivakumar²

^{1,2}Department of Information Technology, Mookambigai College of Engineering, Tamilnadu, India

Abstract: This paper presents a comparison report of two different processes of retrieving a keyword or data's from a given database or from a multiple databases. The process1 known as Extended Boolean Retrieval (EBR) model, it gives us an output from the database. Since EBR model implementation aspects lead to a high cost, we consider an p-norm approach to the EBR implementation. P-norm approach plays a role in the EBR model to maintain strictness of the conjunctions and disjunctions to set them with their own identification on the considerable node. The process2 known as Ensemble Learning Paradigm (ELP). In this paradigm of text categorization aspect, first it assigns a value to a given keyword or data and then starts it's searching process from an index. This value contains the factors such as a position and appearances of word. In existing, they use these concepts in Bag-of-word approach. In this paper EBR model gives an advantage of reformulation aspect, which gives a hundreds or thousands of answers for the given query. In ELP, term frequency identification paves the way to produce a result based on the frequencies of an regarded query in the database. To end, we evaluate with the reported results of these models on query to prove an better retrieving process based on their efficiency and accuracy with the max-score ranking algorithm.

Keywords: Ensemble learning, Index, max-score, rank, term frequency.

I. Introduction

Allowing users to easily search and retrieve a data from database known as retrieval. Though there are many search process to give a better search result, their efficiency and accuracy is unknown to user. The user also unknown of schema of a structured data that has been given by the engine to a user from a server. In such cases, the complex information in a database needs much more descriptions to match their originality to a given query. The properties of the reported results can be failed to be understood by the user. Especially in a domain of biomedical, legal applications, their aspects fails to prove a better efficiency in a result that we received, if it has been then their computational cost leads to infinite level to our extent.

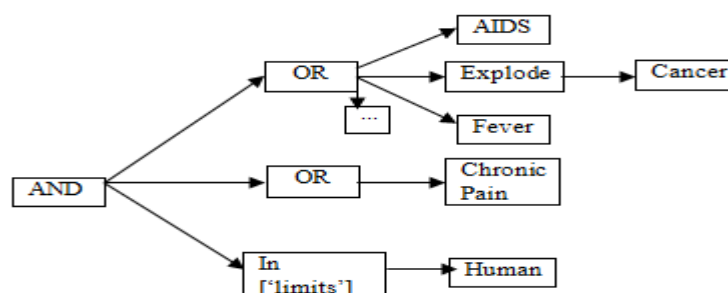


Figure 1: Query tree showing operators functions

Fig. 1 shows how the operators function the retrieval process of Boolean concepts. Here, we have shown the AND-OR format representation for the aspect of efficient extended Boolean retrieval model for the given query.

Thus, we here consider an two better processes of retrieving aspects to prove their efficiency and accuracy. We can choose the domain of biomedical, legal applications for a better result. In these applications, their work of assigning a rank to a database differs from one another and based on their rank basis, the query results have been overviewed and viewed on the output screen. These processes overcome the failure of stop ranking the document one or many of highest ranked result are sufficient based on the constraints assumption.

EBR models are said to be a combination of both the Standard Boolean model and vector space models. Due to the disadvantages in a Standard Boolean model, we go for an EBR model. EBR model has an advantage of weighting factor additionally included into it. The comparison results of Boolean model and EBR model are shown in a test conditions of CSI, CACM and INSPEC. With these comparison bases, we compare an EBR and ELP on their obtained result in conversion to the rank. Here, in an EBR model, we use an p-norm

approach for an implementation basis to reduce an cost. In EBR model, the and-or format representation are considered as an operators to prove their Boolean condition in a searching aspect. Suppose if the user query is $Q=[Aids,Cancer,Fever]$ (1)

Then, the searches are based on an aspect of [Aids AND Cancer AND Fever], [Aids AND Cancer AND NOT Fever], [Cancer OR Fever], etc. Further, the identification of research has noted the query expansion with the extended Boolean retrieval query processing. EBR combines with the content based navigation to give a better result. The appearance of word in a pattern is disregarded because of complexity. These patterns using frequency of a word easier to rank is followed in an EBR. With the p-norm approach implementation of EBR models done with the low cost.

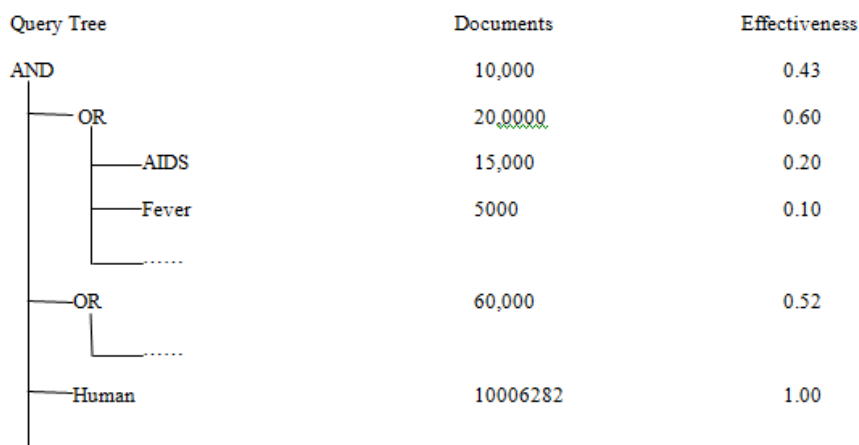


Figure 2: Case study of query given in biomedical application

Fig. 2 shows an example case study by how which the given query in an biomedical application works for selecting a specific disease relevant to the given query from a variety of diseases stored in an single database or multiple database with the patient name exactly relevant to the query that has been given by the user to the parser. Here, we consider an biomedical application as an aspect to undergo our two processes on them with the basic structure to be applied on to them. The mentioned case study result shows the effectiveness of the given keyword or sentence in a searched database and their hierarchy are also to be noted here and their variance activities are also to be noted over here in the biomedical application. In this referred application, the patient name with their regarded complete details are to be extracted along with the search engine that we had proposed to do their ranking based on their identification by using the max-score ranking algorithm in to the EBR and ELP to do their ranking effectively and suggests the top search for the searching keyword or an sentence to a search box.

Ensemble Learning paradigm, which is the first step for the query to execute its access for finding an result, followed by an indexing to note their database format in an order that they need to find an required result for a corresponding query that has been given as input. With the reference of an index, we frame the frequency of a keyword that it occurs in the document and their document frequency is also to be noted for a paradigm that we consider. It overcome response time resulted in Bag-of-word approach and their accuracy in frequency identification. Thus, the resultant of both EBR and ELP are passed through a max-score ranking algorithm for an purpose to assign a priority on the basis of rank obtained from a max-score ranking algorithm. After the rank allocation has been completed the EBR and ELP results are to be compared for a better finding of search prospect.

II. Related work

Searching documents or keywords involved in case of biomedical, legal applications literature. J.H.Lee [1] has developed the different kinds of models for an aspect of retrieval process that the Boolean format permits to flow are proposed. EBR model includes the Boolean queries into it, that has been deeply embedded in the considered process for a domain of biomedical as mentioned in S.Pohl et al. [2]. It includes a conclusion of that Boolean search is not sufficient to be considered, since the binary aspect of matching is not enough because of an indexing errors. Few papers like Karmietal [6] suggests that loosening the strictness consideration of the given query combines with the rank basis priority that we allocate based on a some of the ranking algorithms.

L.Zhang et al. [4], done their research into the performance matrix that they have taken for allocation, such as efficiency, accuracy of the text retrieved from the systems has continued on the inverted index in the sequence basis. Here, they use an method of term-at-a-time methods. As they have preferred an efficient method for preparation, the result still remains sparse as it is. However they produce a result on testing with GOV2 corpus with 61% faster than document time baseline. The Cohen et al. [7] estimated the advantages of an

approach and where they are useful based on classification criteria in a Boolean retrieval result set on the screen that we expect. To practice a classifier, they have taken merely 50% of given documents for judgment. Practically, this leads to an improvement. If these concepts were applied for filtering the recently published documents with their relevance in consideration to systematic reviews more number of documents are to be found. M.E.Smith [8] proposed the aspects of p-norm approach of information retrieval with syntactic query generation, efficiency and theoretical properties as the p-norm model of G.Salton et al. [9] has considered the two characteristics that if absence found, are prepared to retrieval effectiveness. EBR gives a rank basis output from a given query specifications. For the case of ranking in EBR models we have suggested many approaches, one of which is a fuzzy-set models proposed by T.Radecki [5] in case of a document retrieval to include non-binary keyword weights, but considers those effectiveness in the pure Boolean model. There are some of the other approaches implemented in EBR include Waller-Kraft [14], Paice [16], infinite-one [8], inference networks [17], etc.

In an ELP, the process includes various finding query classifications on to it. The R.Bekkerman et al. [10] proposes distributional word clusters verses words for text categorization, which applied to the domain of textual. It results the problem of assuming content of text to the already defined properties. As the value of text consideration improves rapidly on-line and in the corporate domains, which acts as a way to combine the content of text, which paves way as interesting not only from academic but also from a look by the industrial aspects. The M.F.Caropreso et al. [11] gives a learner independent evaluation of the usefulness of statistical phrases for documented text categorization, which faces main problem as training a computer to automatically classify text document, since text categorization has become one of the key aspects for handling and organizing text data.

T.Joachins [12] proposes text categorization with support vector machines (SVM) with many relevant features. It allows achieving high performance. It outperforms strong algorithm word-based setup, which is to be one of the best reported categorization. With these approaches and implementation, their comparison reports are to be overviewed.

III. Our proposal

3.1 Architecture

The user A gives a query to the parser. The parser then continues its action to the database or an server where the data's are stored. After reaching its destination aspects EBR works starts that is Extended Boolean Retrieval model works based on its AND-OR format representation to the given query by using following equations.

$$K_{OR}(N_1, \dots, N_n) = (1/n \sum Q_i^p)^{1/p} \tag{2}$$

$$K_{AND}(N_1, \dots, N_n) = 1 - (1/n \sum (1 - Q_i)^p)^{1/p} \tag{3}$$

The mentioned equations gives us the way by how which the AND-OR representation can be calculated. The two process has to be calculated based on our proposal. It matches and finds by their relevance to keywords and concepts in the given query. It also checks the overall lists popularity. In the optimization activities it identifies whether or not it is being approached for excessive search engine optimization. It looks high up on the page, headings, sub headings, boldface, title, URL format, contents description, ALT tags in case of an graphical format, generic words meta tags, link the given text for inbound lists, etc. By the mentioned categories EBR searches the relevant data on the database. After the searching aspect has been completed in the process1 known as EBR, the reported result has been passed through the max-score ranking algorithm for the purpose of assigning rank to the relevant data's to the given query based on their priority formed by and-or representation to display the regarded result on screen. With this obtained result from process1, system waits for process2 to produce a comparison report.

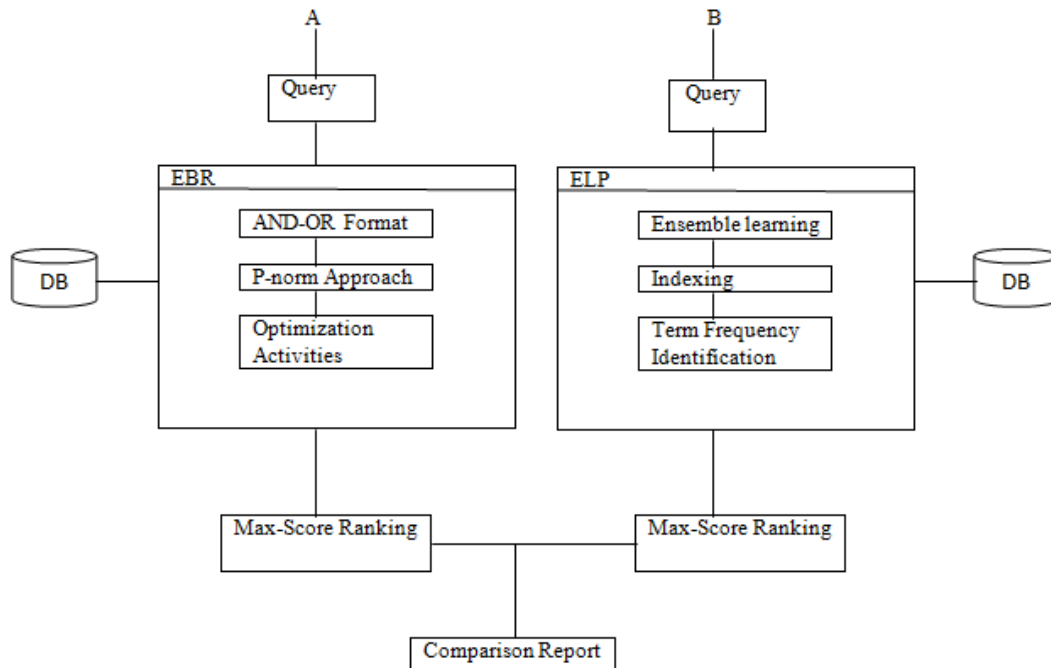


Figure 3: System architecture

Fig. 3 shows the overview of overall system architecture by how which system works to produce a better result for retrieval process.

In the process2 aspects as similar in process1 first the query has to be given by the user B to the parser. The parser then starts its action with Ensemble Learning Paradigm (ELP). The ELP first passes the query to the ensemble learning technique with which the query identification are to be done, what type of query it is given. Then with the completion of the technique, index has to be framed from the database by which the identification or relevance to given query has been noted and used by the ELP for the aspect of searching. The goal of ELP is to figure out during querying process, how many of the relevant documents have been retrieved based on the strictness of their operation. To reduce the complexity in the given query the ELP process transfers the text query to some other easier format to handle can be a vector which is helpful in describing the content of the document. With the completion of index, the next step is to be a term identification by following equations and passed through max-score ranking to assign a rank and this result have been used for comparison with the process1 to produce a result.

$$Tf_{a,b} = n_{a,b} / \sum P_{p,j}^n \tag{4}$$

$$Df_a = \log |D_1| / |d:t_a \in d| \tag{5}$$

$$(tf-df)_{a,b} = tf_{a,b} * df_a \tag{6}$$

3.2 Max-Score Ranking Algorithm

Max-score ranking algorithm is used here for identifying a rank from a hundred or thousands of documents from a database by using a following pseudo code in both the cases of processes ERP and ELP.

```

Maxscore({Aa1, Ab1, ..... Azn}, P)
Initialize min score S and max P
RankLists ← { Aa1, Ab1, ..... Azn }
Forall Aa do
S.push(Ak.curposting(), Ak.curposting().docID);
Ak.next();
While S.isEmpty() ≠ true do
(scoref, docID) ← score(Azn);
p.push(docID, score);
p.pop();
β ← P.pop().score;
update(β, RankLists, S);
    
```

With this referred algorithm, the rank to the following keyword works based on the processes reported the result to the concerned system. As doing this, the example has been considered as shown in the table 1 in case of a biomedical application to search a patient document by using their name as a keyword in the search box.

Table 1: Executed example result

Keyword	Rank	Postings
ANU	10	(D ₁ :2), (D ₂ :10), (D ₄ :2)
RAJA	5	(D ₁ :1), (D ₁ :5), (D ₁₀ :1)....

IV. Retrieval system and control flow

The overall process of the system and their control flow diagram of the retrieval system is shown in fig. 4 implemented by using jdk 1.6. Flow diagram is a graphical representation of the “flow” of data through an information system, modeling its process aspects. Often they are a preliminary step used to create an overview of the system which can later be elaborated. It is also useful for visualization of data processing. Input query has been given by the user A and user B. From the given query two processes are to be followed ERP and ELP. The reported results from ERP and ELP are passed through the max-score ranking algorithm by which the reported results are displayed on screen with their allocated rank. With this ranking obtained from each processes comparison report are to be defined with their performance metrics that they have resulted in the aspect of matching the given query with the relevant data from the database.

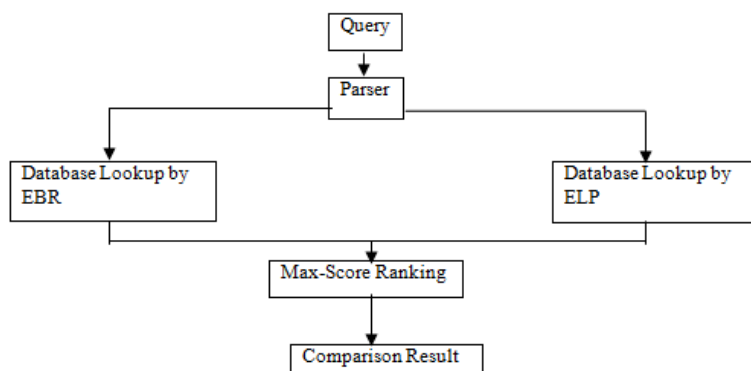


Figure 4: Control flow diagram

V. Performance analysis

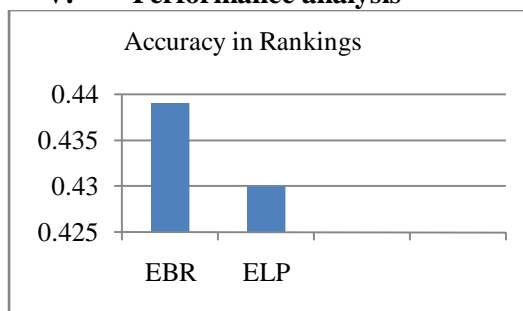


Figure 5: Performance Evaluation in accuracy

Thus, the use of max-score ranking algorithm achieves a time and space exponential gain. Efficiency improvement is obtained in EBR because the data from the database are undergone a Boolean retrieval process compared to ELP which undergone a term frequency style equations to to solve a query. Thus, the graph shows the difference in accuracy of the executed result on variance of two different queries executed in two different process and undergone a ranking priority to the result based on an max-score ranking. In the given graph, the x-axis takes the time taken and y-axis takes the accuracy in rankings. With these aspects, ERB model achieves the better retrieval process when compared with the ELP.

VI. Conclusion

In this paper, we present a system of an comparison report between two retrieval processes such as EBR and ELP models in a biomedical domain. The system can perform static gesture training in comparison of any two retrieval processes. The executed results show that the proposed system allows fast training in finding of different queries at a time with the better processes. The average recognition accuracy of system in allocating a rank to relevant keywords from database is 79.8%. Term weight has an advantage that the document can itself calculate its score, which is of similar. The system gives the result that ERB model achieves the better retrieval

process when compared with the ELP. The only limitation of our system is that, the response time takes longer due to the conversion aspect, search allocation based on priority and rank assumption. Future work will include extending the developed system with other different ranking algorithms.

References

- [1] J. H. Lee, "Analyzing the effectiveness of extended Boolean models in Information Retrieval," *Cornes University, Tech. Rep. TR95-1501*, 1995.
- [2] S. Pohl, J. Zobel, and A. Moffat, "Efficient Extended Boolean Retrieval," *University of Melbourne*, 2012.
- [3] S. Pohl, J. Zobel, and A. Moffat, "Extended Boolean Retrieval for systematic biomedical reviews," in *proc. of the 33rd Australian Computer Science Conf. (ACSC 2010), ser. Conf. in Research and Practice in Information Technology (CRPIT), vol. 102*. Brisbane, QLD, Australia: Australian Computer Society, Jan. 2010.
- [4] L. Zhang, I. Ajiferuke, and M. Sampson, "Optimizing search strategies to identify randomized controlled trials in MEDLINE," *BMC Med. Res. Meth.*, vol. 6, no. 1, p. 23, May 2006.
- [5] T. Radecki, "Fuzzy set theoretical approach to document retrieval," *Inform. Process. Manag.*, vol. 15, no. 5, pp. 247-259, 1979.
- [6] S. Karimi, S. Pohl, J. Zobel, and F. Scholer, "The challenge of high recall in biomedical systematic search," in *proc. of the 3rd Int. Workshop on Data and Text Mining in Bioinformatics. Hong Kong, China: ACM, Nov. 2009*, pp. 89-92.
- [7] A. M. Cohen, W. R. Hersh, K. Peterson, and P. Y. Yen, "Reducing workload in systematic review preparation using automated citation classification," *J. Am. Med. Inform. Assoc.*, vol. 13, no. 2, pp. 206-219, 2006.
- [8] M. E. Smith, "Aspects of the p-norm model of information retrieval: Syntactic query generation, efficiency, and theoretical properties," *Ph.D. dissertation, Cornell University, May 1990*.
- [9] G. Salton, E. A. Fox, and H. Wu, "Extended Boolean Information Retrieval," *Commun. ACM*, vol. 26, no. 11, pp. 1022-1036, Nov. 1983.
- [10] R. Bekkerman, R. El-Yaniv, N. Tishby, and Y. Winter, "Distributional Word Clusters versus Words for Text Categorization," *J. Machine Learning Research*, vol. 3, pp. 1182-1208, 2003.
- [11] M. F. Caropreso, S. Matwin, and F. Sebastini, "A Learner-Independent Evaluation of the Usefulness of Statistical Phrases for Automated Text Categorization," *Text Databases and Document Management: Theory and Practice*, A.G. Chin, ed., pp. 78-102, Idea Group Publishing, 2001.
- [12] T. Joachims, "Text Categorization with Support Vector Machines: Learning with Many Relevant Features," *Proc. 10th European Cong. Machine Learning (ECML '98)*, pp. 137-142, 1998.
- [13] J. H. Lee, "Properties of extended Boolean models in Information Retrieval," *Proc. 17th Annual International ACM SIGIR Conf. in Research and Development in Information Retrieval*, pp. 182-190, 1994.
- [14] W. G. Waller and D. H. Kraft, "A mathematical model of a weighted Boolean retrieval system," *Inform. Process. Manag.*, vol. 15, no. 5, pp. 235-245, 1979.
- [15] William Hersh, *Information Retrieval: A Health and Biomedical Perspective*, Springer, 3rd edition, Nov 2008.
- [16] C. D. Paice, "Soft evaluation of Boolean search queries in Information Retrieval systems," *Inf. Technol. Res. Dev. Appl.*, vol. 3, no. 1, pp. 33-41, Jan 1984.
- [17] H. Turtle and W. B. Croft, "Inference networks for document retrieval," in *Proc. of the 13th Ann. Int. ACM SIGIR Conf. on Research and Development in Information Retrieval. Brussels, Belgium: ACM, 1990*, pp. 1-24.