# Priority Based Retrieval of Information with Documents Domain Division Approach

**Alla Hashwanth, P.V.N Sanghavi, B.B.V Satya Vara Prasad, Ravi Kumar Tenali**

*Abstract: Human beings "search" for accomplishing their objectives (desires) in all sorts of situations. The research is still continuing to gain more quality in searching in order to more easily acquire require data and information in all aspects. With most of the existing techniques, the irrelevant documents in searching process are also retrieved because of improper application of knowledge about corresponding domain to which the documents searching belongs. In this paper, an efficient method using which the documents submitted to the database will be divided into specific domain sis proposed. When the users searches for a particular domain or documents, the keywords used in the search used for the identification of domain, and the relevant documents are only retrieved on priority basis. The searched keywords along with association rules, decision trees, indexing techniques plays major role in building this new method. The documents which have more count value of keyword will be retrieved first and later documents in accordance with the descending order of count value of keywords present in them, which results in the retrieval of more related documents at the top.*

*Index Terms: Keyword search, document division, domain specific, Priority, Count value, MySQL, Databases, Information retrieval. Decision tree, Inverted indexing.*

## I. INTRODUCTION

Due to the huge amount of data available in the web, the retrieval of most related information has become very difficult now a days. The information retrieval is more important because the famous search engines like Google, Yahoo, Bing etc., are using this kind of techniques only to provide information to the users or searchers. Users may not know the exact keyword that was enabled in the database. So, construct a technique in which keywords are extracted from user typed keyword in the search bar and finds the reputation of these keyword in the documents stored in the database. According to the standard text mining process, document is indicated as vector of whose dimension is number of the discrete keywords present in it, having the high frequency of keyword count[1].Thereupon standard text classification could be competitively costly.

Sensing the keyword automatically for retrieving the text and Innovative ranking is applied to build vigorous methodology and upgrades Information retrieval system[2]. D.Tumeretal[3] has decided an exact assessment on Semantic methodology and inquiry execution of Keyword and Semantic based Search Engines like Yahoo, Google, Hakia and Msn and so on. Execution of data recovery frameworks depends on their exactness proportion and characteristic language queries[4]. As new things show up, the classification is increasing day by day , so our application should meet this classification to retrieve the accurate information. Data Retrieval (IR)[5] is discovering material of an unlabeled nature(usually message) that fulfills a data need from inside huge accumulations. In order to parse and label an archive string ( parses a content record into tokens) change the tokens to lexemes [6]. Structure of an insightful specialist for substance based ordering and recovery of applicable archives from a substantial gathering, for example, the web [7], M Andago [8], chosen a assessment of the Semantic Search Engine and contrasted a Keyword Search Engine by utilizing the equation "initial 20 Precision". He has picked thirty questions and gone into the web crawlers and determined precision proportions. The Google outflanks with Hakia. Because it has precision of 0.64 higher compared with Hakia. So, it was clarity that Google is has higher precision than Hakia. Existing classification techniques, are mostly based on the search engines, and these search engines are keyword based mostly. In different words, every document consists of a group of meaningful terms (also referred to as descriptors or keywords) that area unit believed to specific to the content present in document. These keywords area unit assigned some weights depending on factors like their frequency of keywords, incidence (i.e. using Boolean vector based, or probabilistic methods; see [9],[10],[11]). Keyword based search engine is very helpful for finding the information present in the internet. It cannot find the meaning of some terms and expressions which are used in the webpages. Currently keyword based search approach has reached a plateau. According to the literature surveys 25% of searching in the web do not give the accurate results because of due to the increase of sixty-terabyte in the size of the web daily[12]. This paper focuses on domain specific division of documents which is not much focused in survey. The division is done by using decision tree algorithm which is easy to increase more accuracy than those of by using support vector machines(SVM's)[14] as previous research is done by using (SVM's).

*Retrieval Number: F2553037619 /19©BEIESP*
*Journal Website: www.ijrte.org*

562

*Published By:*
*Blue Eyes Intelligence Engineering*
*& Sciences Publication*

## II.   METHODOLOGY

**PROPOSED SYSTEM DESIGN**

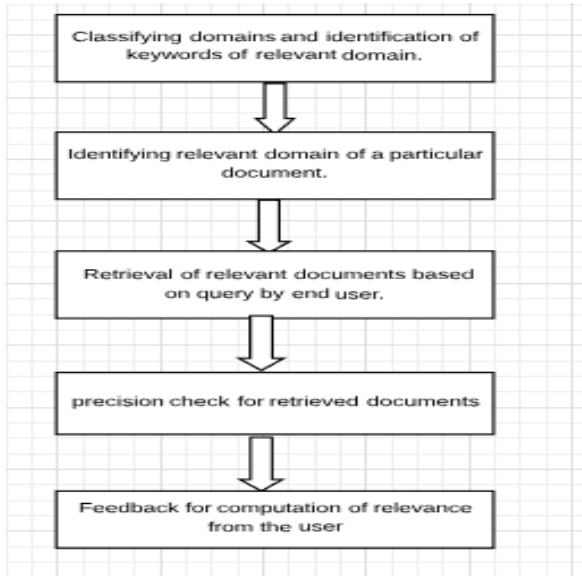The work will be carried out as follows:



**Fig.1: Proposed system algorithm**

### II.1 Classifying domains and identification of keywords of   relevant domain:

Domain classification using predetermined keywords is the process in which division of the domains in the data base was done by storing the expected keywords for each domain for example : computer science consists of keywords like computer, algorithm, data structures, big 'O' notations etc., Similarly many domains like networking, signals, software engineering etc., are created. But when the user enters some document which does not contain any keywords that was present in our database so, in such cases. Bootstrap provides learning methodology for the machine to generate new keywords from the user submitted documents and are mapped with more relevant domains.

A. *Working methodology of bootstrap:*

Bootstrapping is an iterative process in which less amount of data is given for learning then the system will develop by using bootstrapping which gives more accuracy output.

At first we will give input as less number of keywords then from the unlabeled data at each iteration a label data is outputted which matches this keywords. After the completion of this iterations a new  model  is developed from those labeled outputted data and this new model will identify the new labels from the unlabeled data and this process will repeat for every unlabeled data to be  labeled

So, by the above process with less amount of input data set we can  extract  many  labeled data sets which will be taken as decision identification keywords by the decision tree as a result more accuracy will be increased in the output.

For the lion consists of the activities like roaring, hunting etc., By using this activities we must identify the label data

as lion or the output as lion, this is how bootstrapping works.

### *Identifying relevant domain of a particular document:*

After the division is done, when the user submits his/her documents then the decision tree algorithm will identify to which domain the submitted document belongs to by sensing those predetermined keywords from the document , If there are more relevant keywords of computer science domain specific in a document then it will be stored in the computer science domain in the database, similarly the domain identification for the other documents also done and stored in specific domain tables in the database.   In some cases may have equal frequency of keywords that are relevant to more than one domain in a document. In such cases, the document will be stored in multiple domains.

A. *Working principle of decision tree:*

The decision tree is nothing but a decision maker for example an event can happen or  not by predicting whether a data consists of particular word or not. Here each leaf stump of the tree is a if condition which means a decision taker .So, each leaf stump is connected to different type of domain . For, example a decision leaf stump which filters the computer domain specific documents will checks whether the document really consists of those computer specific words like computer, computer science, programming etc,. So, that the filtered document will be sent to that computer specific domain .In this way, the decision tree will checks many documents to classify those documents  into their specific domains.

If condition1 and condition2 and condition3 then outcome.

Among the available decision making tools decision tree technique has many benefits such as:

1. Decision making trees are easy to understand by the normal people after a brief explanation to them about decision making trees.

2. Decision tree has better way of its alternatives, the possibility of giving the accurate results was more , and costs of its implementation are less as compared with other techniques

3. It is used to get  the best, worst and expected values for different situations.

4. Decision tree can be combined very easily with the other decision making techniques.
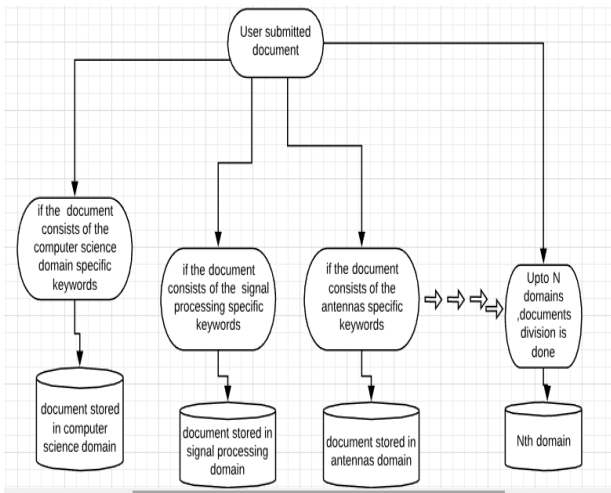
**Fig.2: Decision making tree**

## II.2 Retrieval of relevant documents based on query by end user

When user types the sentence to search a document, then the keywords will be extracted from that searched sentence by avoiding the phrases, and will identify to which domain the keywords belongs to. After knowing the keywords domain, the decision tree is used to retrieve documents from that particular domain and the retrieving will be done in priority wise based on the density of keywords present in the each document and the density of those keywords in each document can be obtained by using inverted indexing technique.

To extract the keywords from the user searched sentence inverted index technique is used as follows:
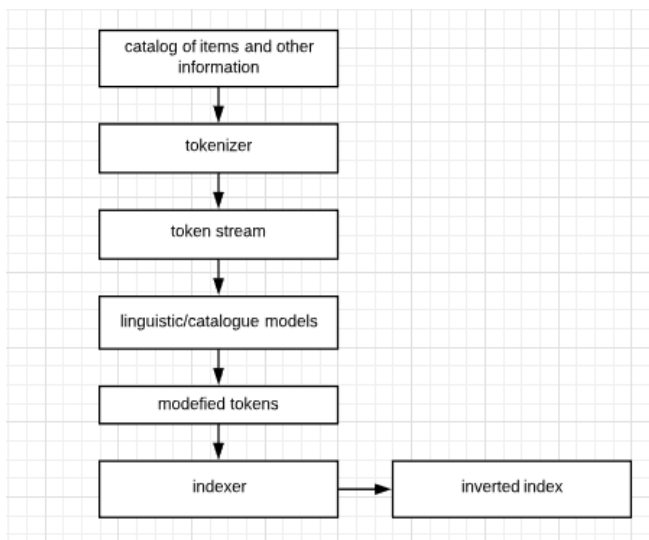


**Fig.3: Algorithm for inverted indexing.**

A. *Explanation about the working of inverted index:*

At first the user searched query or information are sent to tokenizer, this tokenizer will makes this information or break downs the information into token streams for example user query is > what is the definition of definition of computer science then this sentences will be tokenized into tokens as follows What|| is|| the|| definition|| of|| computer|| science then the linguistic catalogue models will

avoid all the phrasal verbs like of, the, this etc., and extract the major keywords like computer, science, definition etc,. then this modified tokens are sent to indexer , then this indexer will identify how many times the each keyword was repeated in different documents and then the documents which has the most repetition of this keywords are retrieved first later on the descending order. And resulted output is in the form of inverted index. As we know that the inverted indexing gives us priority wise retrieval of information according to the keyword frequency.

### II.3 Feedback from the user:

Feedback form is provided for both the document submitting and document retrieving users as follows:

**Feedback form for submitted users:** Displaying of to which domain the user submitted document belongs to, if the domain is falsely identified by the machine then user is asked to suggest to which domain his/her document belongs and also will be asked for some suggested keywords from the user. This feedback from the user will be used in improvement of proper identification of the domain.

**Feedback from for searched users:** Here, after the retrieving of the documents based on the user typed content in the search bar for searching. Then takes the feedback from the user whether his/her required document is retrieved or not, if he did not receive it then feedback is asked from him to suggest which domain specific document is he searching for, based on that feedback our search results will be improved for further searching users.

### II.5 Precision Check

In order to check our results with the existing models precision and recall techniques are used as follows:
Precision: - Precision is the ratio of Documents with respect to total relevant documents retrieved to the irrelevant documents.
Recall: - Recall is the ratio of documents with respect to retrieved relevant document to the possible relevant documents.

### III. RESULTS

To improve the accuracy in the resulted output of decision tree, implementation of bagging for decision tree technique is done. Bagging: By using the bagging technique we can increase the accuracy of our decision tree because, as the number of base line nodes increases then the number of decision makings which will filters out the output will increases.As the below table and graph shows you the practical implementation which increases the accuracy according to the increase in number of nodes or leaf stump number.Here's however the sacking ensemble performs once varied
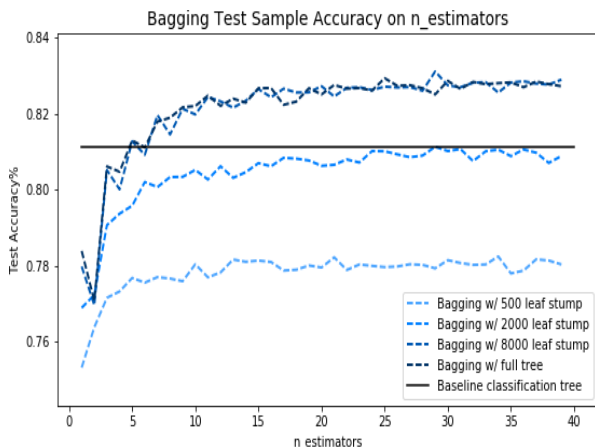
(1) the quantity of base estimators and (2) the scale of the bottom estimators (no. of leaves)

Results:

| Leaf stump number | Accuracy at max Estimator | Maximum number of estimators |
|---|---|---|
| 500 | 0.78 | 40 |
| 2000 | 0.81 | 40 |
| 8000 | 0.83 | 40 |
| Full tree | 0.834 | 40 |
| Base line tree | 0.81 | 40 |



Precision for our system is

Precision = (retrieved documents /estimated

documents)*10

=17/20=0.85*100 = 85% accuracy  .

## II CONCLUSION

The results effort proves that evaluation parameters like accuracy in decision tree, precision and recall plays a vital role in the improvement of accuracy in both domain division of documents and giving accurate results for searching user. Priority wise retrieval will benefits the search user as his more accurate resulted information is retrieved at the top based on learning and resulting techniques of bootstrapping. Hybrid system will mostly increases the performance of system . So in the next research scope is to implement this hybrid methodology to increase more and more the accuracy over than  95% in obtaining the output so that it satisfies the user.

## FUTURESCOPE

This implementation is very useful in medical, because when the patient submit some symptoms of  his disease, it will automatically identify the with what disease the patient is suffering from and it can also retrieve the suggested medicines required to cure his disease. The image files domain identification is very useful in to identify the disease of a person by simply scanning the body skin from his photo.

This is also very useful in classification of political problems faced by the common people asfollows:

When common people submitted their problems or issues to their political leaders through online, this machine will directly identifies the problem in the document which reduces the reading time of that document by the political leader which results in quick decision making by them.

**REFERENCES**:

1. Shashank Pandit, Soumen Chakrabarti, Varun Kacholia,
2. S.Sudarshan "Bidirectional Expansion For Keyword Search on Graph Database", Proceeding of 31st VLDB Conference 2005.
3. YannisPapakonstantinou, VagelisHristidis, Luis Gravano, "Efficient IR-Style Keyword Search over Relational Databases", Proceedings of the 29th VLDB Conference, Berlin, Germany,2003.
4. M.A.Shah, D.Tumer and Y. Bitirim "An Empirical Evaluation on Semantic Search Performance of Keyword- Based and Semantic Search Engines: Google, Yahoo, Msn and Hakia," Fourth International Conference on Internet Monitoring and Protection, 2009.
5. D.K.Lobiyal, .V Mala," Semantic and Keyword Based Web Techniques in Information Retrieval", International Conference on Computing, Communication and Automation(ICCCA2016).
6. http://web.stanford.edu/class/cs276/ (active May2016).
7. https://www.compose.io/articles/indexing-for-full-text-searchin-postg resql/ (active May2016).
8. F. blarir, F. Meziane, N. K. Mimouni,  "An Intelligent Agent for Content-Based Indexing and Retrieval of Documents", Fourth International Conference a knowledge- Based Intelligent Engineering Systems and Allied Technologies, Wh Aug-1" Sept 2000,Brighton.
9. P.L Phoebe and A. M Thanoun.M. Andago,, "Evaluation of a Semantic Search Engine against a Keyword Search Engine Using First20 Precision," In Proceedings of the 29th Annual International Conference on Research and Development in Information Retrieval, ACM Press,pp.735–746, 2010.
10. G.Z, Liu, "Semantic vector space model: implementation and evaluation." Journal of the American Society for Information Science, 48(5), 395-417,1997.
11. R., Korfhage" Information Storage and Retrieval. "John Wiley and Sons, London,1997.
12. R.M, Comparing "Boolean and probabilistic information retrieval systems across queries and disciplines," JASIS, 48(2), 143-156,1997.
13. Prof Sukhjit Singh Sehra, Bhoomika, Prof Anand Nayyar " A Review paper on Algorithms used for Text Classification" International Journal of Applicationor Innovation in Engineering & Management (IJAIEM) Web Site: www.ijaiem.org Email: editor@ijaiem.org,editorijaiem@gmail.com Volume 2, Issue 3, March 2013.
14. K P.Ukey, Dr. A.S. Alvi " Text Classification using Support Vector Machine", International Journal of Engineering Research & Technology (IJERT) Vol. 1 Issue 3, May – 2012.

**AUTHORS PROFILE**

AllabHashwanth is a UG student in koneru lakshmaiah Education Foundation and belongs to the branch Electronics and Computer Engineering. He is interested in the domain of Web development and Text mining.

P.V.N.Sanghavi is a UG student in koneru lakshmaiah Education Foundation and belongs to the branch Electronics and Computer Engineering. She is interested in the domain of Web development and Text mining.

*Retrieval Number: F2553037619 /19©BEIESP*
*Journal Website: www.ijrte.org*

565

*Published By:*
*Blue Eyes Intelligence Engineering*
*& Sciences Publication*

MR. B.B.V.SatyaVara Prasad is working as an Assistant Professor in Department of Electronics and Computer Engineering in Koneru Lakshmaiah Education Foundation. He presented the several research papers in reputed international journals and he attended several national and international conferences. His area of interest is Computer Networks and Data mining.

An efficient Assistant Professor, received M.Tech (C.S.E) from Swarnandra College of Engineering and Technology (JNTUK) .working as an Assistant Professor in Department of ECM, Koneru Lakshmaiah Education Foundation (KLEF) .He has 14 years of teaching experience. He has published many papers in International Journals & his areas of Interest includes Computer Networks, Data Mining and Cloud computing.

*Retrieval Number: F2553037619 /19©BEIESP*
*Journal Website: www.ijrte.org*

566

*Published By:*
*Blue Eyes Intelligence Engineering*
*& Sciences Publication*