

Comparative Analysis of Intersection Algorithms on Queries using Precision, Recall and F-Score

Sachin Kumar

Sr.Asst.Prof – Computer Science/IT
Fairfield Institute of Management & Technology
(Affiliated to GGSIP University, New Delhi)

Pratishtha Gupta, PhD

Asstt. Professor – Department of Computer Science
Banasthali Vidyapith,Jaipur,India

ABSTRACT

Search Engines are the software systems to search for the required information on the World Wide Web. Though Search Engine Optimization deals with increasing the visibility of Web pages on search results, this research work focuses on the search algorithm, which is expected to bring the most appropriate, relevant and required documents for the given search query. The demand of the new such algorithm has been presented by a comparative study of intersection algorithm and its variants. CRANTOP database is a standard database available on WWW, which not only provides the database of documents and queries but also provide the actual relevance of documents for the corresponding queries. Study has been done on a sample part of that database, which is available in appendix and analysis has been done on the basis of Precision, Recall and F-Measure values.

Keywords

Precision, Recall, F-measure, Algorithms, Database et al.

1. INTRODUCTION

Algorithm used for Information Retrieval impacts the customer satisfaction a lot. Many a times, lot much data required by the user is available but cannot be retrieved just because the algorithm used could not understand the requirement of the user. This satisfaction of the user requirements is measured in many ways. Precision, Recall and F-Score are some of those measurement models, through which effectiveness of an algorithm can be measured.

This paper emphasizes on the evaluation of a query from a standard database CRANTOP, using a simple intersection algorithm and its variants. Comparisons in the results have been done to emphasize the fact that differences exist in the outcomes of the same query on the same database.

2. REVIEW OF LITERATURE

According to a study conducted by [3] included theoretical, methodological and empirical aspects to explore modalities of use and suggest further avenues. However [4] feels that the objective of this paper is to characterize the changes in the rankings of the top ten results of major search engines over time and to compare the rankings between these engines. [2] suggested in his studies that term position information, as provided in some Boolean systems in the form of field restriction and term proximity, is reviewed and its value assessed. Non-Boolean retrieval in the form of the ranked output experiment has not so far used term position information but has concentrated on schemes of term weighting.

[1] revealed that in the first approach we set the same threshold for each of the IF systems. In the second approach the threshold of each IF system is tuned independently to

maximise its own EP (“local optimisation”). In the third approach the thresholds of the IF systems are jointly tuned to maximise the EP of the combined system (“global optimisation”). [7] found that an interdisciplinary framework for list, especially information retrieval (IR), in a way that goes beyond the cognitivist ‘information processing paradigm’. The main problem of this paradigm is that its concept of information and language does not deal in a systematic way with how social and cultural dynamics set the contexts that determine the meaning of those signs and words that are the basic tools for the organisation and retrieving of documents in List.

[6] founded that the objective of the paper is to amalgamate theories of text retrieval from various research traditions into a cognitive theory for information retrieval interaction. Set in a cognitive framework, the paper outlines the concept of polyrepresentation applied to both the user's cognitive space and the information space of IR systems. The concept seeks to represent the current user's information need, problem state, and domain work task or interest in a structure of causality. [8] said that the purpose of this article is to review the research on human-computer interfaces for library-based commercial online information retrieval (IR) systems.[10] suggested that when developing an IR system it is necessary to think of it as a controllable system. It makes the system capable of fulfilling this system. It leads to its requirements for an optimal search and optimization is one of the problem of control.

[9] feels that the purpose of this paper was to study users' behaviour when using different search engine results pages (SERPs) to identify what types of scents (cues) were the most useful to find relevant information to complete tasks on the Web based on information foraging theory. [5] described the design and implementation of a system for computer generation of linked HTML documents to support information retrieval and hypertext applications on the World Wide Web. The approach is based on work by Salton and others, but extends the concept to be compatible with the World Wide Web browser environment by adding an interactive indexing technique that is well suited to the mouse-based point-and-shoot input common to windowed browsers. The system does not require text query input, nor any client or host processing other than hypertext linkage.

The goal of this work is to construct a fully automatic system in which original text documents are read and processed by a computer program that generates HTML files, which can be used immediately by Web browsers to search and retrieve the original documents. Thus, a user with a large collection of information — for instance, newspaper articles — can feed these documents to the program described here and produce directly, without further human intervention, the necessary files to establish World Wide Web home and related pages, to

support interactive retrieval and distribution of the original documents.

2.1 Intersection Algorithm on Information Retrieval

The intersection operation is the crucial one; it needs to efficiently intersect postings lists so as to be able to quickly find documents that contain both terms. There is a simple and effective method of intersecting postings lists using the intersection algorithm. It maintains pointers into every lists and walk through the every postings list simultaneously, in time linear in the total number of postings entries. At each step, it compares the docID pointed to by every pointer. If they are the same, put that docID in the results list, and advance every pointers. Otherwise, advance the pointers pointing to the smaller docID. If the lengths of the postings lists are x and y , the intersection takes $O(x+y)$ operations. Formally, the complexity of querying is $O(N)$, where N is the number of documents in the collection. Its indexing methods gain it just a constant, not a difference in O time complexity compared to a linear scan, but in practice the constant is huge. To use this algorithm, it is crucial that postings be sorted by a single global ordering. Using a numeric sort by docID is one simple way to achieve this. This algorithm can be extended to process more complicated queries like for example.

(Brutus OR Caesar) AND NOT Calpurnia

Query optimization is the process of selecting how to organize the work of answering a query so that the least total amount of work needs to be done by the system. A major element of this for Boolean queries is the order in which postings lists are accessed. For each of the t terms, it needs to get its postings, then AND them together. The standard heuristic is to process terms in order of increasing document frequency, if it start by intersecting the two smallest postings lists, then all intermediate results must be no bigger than the smallest postings list, and we are therefore likely to do the least amount of total work, Christopher D. M. et al [2005].

It keeps the frequency of terms in the dictionary it allows to make this ordering decision based on in-memory data before accessing any postings list. It will get the frequencies for all terms, and it can then estimate the size of each OR by the sum of frequencies of its disjuncts.

2.2 Intersection via Skip Pointers Algorithm on Information Retrieval

Here also like intersection extensions to postings list data structures and ways to increase the efficiency of using postings lists. It recalls the basic postings list intersection operation if it walks through two or more postings lists simultaneously, in time linear in the total number of postings entries. If the list lengths are m and n , the intersection takes $O(m+n)$ operations. Can it be better than this? That is empirically, it can usually process postings list intersection in sub linear time. It can be possible if the index is not changing too fast, Christopher D. M. et al [2005].

One way to do this is to use a skip list by augmenting postings lists with skip pointers at indexing time. Skip pointers are effectively shortcuts that allow to avoid processing parts of the postings list that will not figure in the search results. The two questions are arises here that where to place skip pointers and how to do efficient merging using skip pointers. Here a number of variant versions of postings list intersection with skip pointers is possible depending on when exactly it checks the skip pointers. Its results is again similar to the intersection

algorithm kind of results as shown below in query 1 to 50 and the recall and precision parameters values are same.

2.3 Positional Intersection Algorithm on Information Retrieval

Here, for each term in the vocabulary, it stores postings of the form docID:(position1,position2,...), where each position is a token index in the document. Each posting will also usually record the term frequency. To process a phrase query, it still needs to access the inverted index entries for each distinct term. Here it finds the results where there are docID on the same position in the document in every postings list in a query that is this query will be retrieved otherwise it cannot be retrieved if the same positions are not found in the postings list of any of the queries in the database, Christopher D. M. et al.[2005].

3. PROPOSED ALGORITHM STEPS

- 1.) Break Query :- Break-row queries are often loosely referred to as terms or words, but it is sometimes important to make a type distinction.
- 2.) Phrase Queries: - A phrase is a collection of words in sentence in a document like “A cakes and bakers manufacturing company has many distributors” is not a match. Most recent search engines support a double quotes syntax (“manufacturing company”) for phrase queries.
- 3.) Two –Words Phrases :- The two words query is the two- words phrases in that one approach to handling phrases is to consider every pair of consecutive terms in a document as a phrase.
- 4.) Three –Words Phrases:- The three words query is the three- words phrases in that it has increased the priority of a root words in the database to be retrieved break them a phrase in single root words so that they can retrieve accurate amount of information from the database.
- 5.) Remove Stopwords :- Common words which would appear to be of little value in helping select documents matching a user need are excluded from the vocabulary entirely. These words are called stop words.
- 6.) Convert Remaining Into Root Words :- Here may be a root words like representation, represented or representing. It will make the accurate data to be retrieved from the set of database as a root word.
- 7.) Classify Into Specific And General Words :- General words and specific words are not opposites instead, they are the different ends of a range of words. General words refer to groups whereas specific words refer to individuals.
- 8.) Synonyms :- Similar meaning root words are called synonyms
- 9.) Intersection With Specific Words :- Intersection have taken as the common postings together as if it have two Boolean queries located in the dictionary. Then intersection operation need to efficiently intersect postings lists so as to be able to quickly find documents that contain both terms.

4. CRANTOP DATABASE

Crantop database is an xml based database in which this research is carried out. It contains 3 parts that is documents with their Doc IDs in their dictionaries. This database would contains Doc.ID 1 to Doc.ID 1400 in which there are several nodes of linked list structure which can be found by the help of softwares like MS-Word. In the second database it contains different queries i.e. the research would be carried out on query 1 to 50 its analysis would be carried out in MS-Excel by sorting the every words into ascending order and make it into a new sparse matrix representation. In the third database its every documents ranking would be given that is how many times a single link list can be carried out in a single query from 1 to 50 for the Doc ID's 1 to 1400 and its rankings would be given here on the basis of that its precision and recall can be calculated here with its f-measure with different Intersection, Intersection with skip pointers and positional intersection algorithms and based on that new algorithm would be proposed and based on that second table would be analysed. This is the need of this crantop database by this research can be carried out. This database is available at the link given below <ftp://ftp.cs.cornell.edu/pub/smart/cran/>

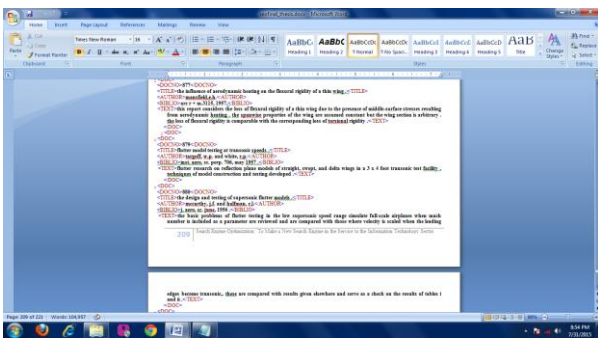


Fig1: Snapshot of Documents of CRANTOP database

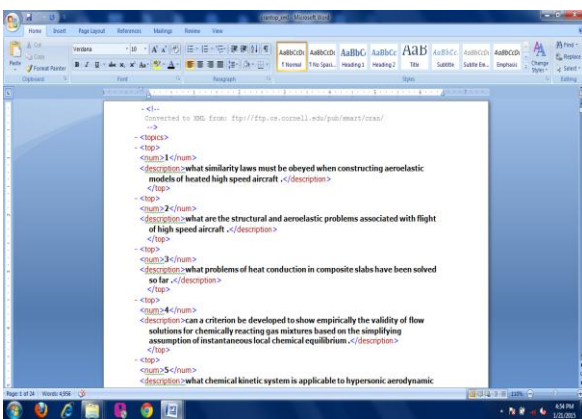


Fig2: Snapshot of Queries of CRANTOP database

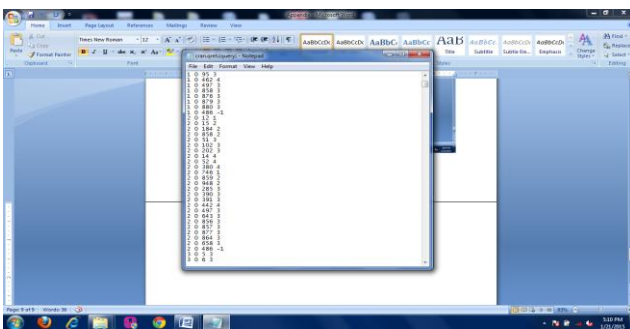


Fig3: Snapshot of relevance result of CRANTOP database

5. PRECISION AND RECALL PARAMETERS

The two most frequent and basic measures for information retrieval effectiveness are precision and recall. Information retrieval system returns a set of documents for a query.

Precision (P) is the fraction of retrieved documents that are relevant

$$\text{Precision} = (\text{relevant items retrieved}) / (\text{retrieved items})$$

Recall (R) is the fraction of relevant documents that are retrieved

$$\text{Recall} = (\text{relevant items retrieved}) / (\text{relevant items})$$

Typically web surfers would like every result on the first page to be relevant (high precision) on the other hand professional searchers such as paralegals and intelligence analysts are very concerned with trying to get as high recall as possible, and will tolerate fairly low precision results in order to get it, Christopher D. M. et al[2005].

A single measure that trades off precision versus recall is the F measure, which is the weighted harmonic mean of precision and recall.

$$\text{F-measure} = 2 * \text{Precision} * \text{Recall} / (\text{Precision} + \text{Recall})$$

The advantage of having the two numbers for precision and recall is that one is more important than the other in many circumstances .Typically, web surfers would like every result on the first page to be relevant (high precision) but have not the slightest interest in knowing let alone looking at every document that is relevant. In contrast, various professionals searchers such as paralegals and intelligence analysts are very concerned with trying to get as high recall as possible, and will tolerate fairly low precision results in order to get it. Recall is a non-decreasing function of the number of documents retrieved whereas precision usually decreases as the number of documents retrieved is increased. Then here a single measure that trades off precision versus recall is the F-measure, which is the weighted harmonic mean of precision and recall Christopher D. M. et al[2005].

To get more high precision and recall we use the proposed algorithm which I will be discussed here step wise and it will show the difference between these three algorithms which are discussed here and the proposed one. The proposed algorithm stepwise to be discussed above.

There is a good reason why we use precision, recall and f-measure to measure the proposed algorithm for information retrieval problems. In almost all circumstances, the data is extremely skewed : normally over 99.9% of the documents are in the nonrelevant category. A system tuned to maximize accuracy can appear to perform well by simply deeming all documents nonrelevant to all queries. Even if the system is quite good trying to label some documents as relevant will almost always lead to a high rate of false positives. Users are always going to want to see some documents and can be assumed to have a certain tolerance for seeing some false positives providing that they get some useful information. The measures of precision and recall here concentrate the evaluation on the return of true positives, asking what percentage of the relevant documents have been found and how many false positives have also been returned. Again, there will be a conclusion that if we measure it on CRANTOP database with intersection algorithms and its alias its precision, recall and f-score are less as compared to the

proposed algorithm will make it is high to all the parameters like precision, recall and f-score the comparisons are given in below two tables.

6. MEASUREMENT OF PRECISION AND RECALL AT START OF THE COMPARATIVE ANALYSIS WITH F-MEASURE

It is analyzing that table on intersection algorithm and its another alias on 1400 documents and 50 queries. They are finding it through CRANTOP database which are relevant documents and how many documents are retrieved from the total relevant documents from the total documents found and now finding the Precision, Recall and F-score which are very low as compared to the proposed algorithm.

Table1. Measurement of Precision & Recall with F-measure

| Q. N O | Intersection algorithm and its alias | | | | Precision | Recall | F-score |
|--------|--------------------------------------|---------------|---------------|-----------------|-----------|--------|---------|
| | Relevant doc | Doc retrieved | Relevant doc. | Total doc found | | | |
| 1 | 0 | 0 | 29 | 29 | 0 | 0 | 0 |
| 2 | 0 | 0 | 25 | 25 | 0 | 0 | 0 |
| 3 | 0 | 0 | 9 | 9 | 0 | 0 | 0 |
| 4 | 0 | 0 | 3 | 3 | 0 | 0 | 0 |
| 5 | 0 | 1 | 4 | 5 | 0 | .25 | 0 |
| 6 | 2 | 11 | 9 | 5 | .182 | 2.2 | 0.34 |
| 7 | 0 | 0 | 6 | 6 | 0 | 0 | 0 |
| 8 | 0 | 1 | 5 | 5 | 0 | .2 | 0 |
| 9 | 0 | 0 | 4 | 4 | 0 | 0 | 0 |
| 10 | 0 | 0 | 9 | 9 | 0 | 0 | 0 |
| 11 | 1 | 0 | 8 | 8 | 0 | .125 | 0 |
| 12 | 1 | 0 | 7 | 7 | 0 | .143 | 0 |
| 13 | 0 | 0 | 5 | 5 | 0 | 0 | 0 |
| 14 | 2 | 8 | 3 | 3 | .25 | 2.67 | 0.46 |
| 15 | 1 | 0 | 3 | 3 | 0 | .333 | 0 |
| 16 | 1 | 0 | 4 | 4 | 0 | .25 | 0 |
| 17 | 1 | 0 | 3 | 3 | 0 | .333 | 0 |
| 18 | 1 | 0 | 4 | 4 | 0 | .25 | 0 |
| 19 | 0 | 0 | 10 | 10 | 0 | 0 | 0 |
| 20 | 0 | 0 | 10 | 10 | 0 | 0 | 0 |
| 21 | 1 | 0 | 5 | 5 | 0 | .2 | 0 |
| 22 | 1 | 0 | 2 | 2 | 0 | .5 | 0 |
| 23 | 0 | 0 | 33 | 33 | 0 | 0 | 0 |
| 24 | 1 | 0 | 4 | 4 | 0 | .25 | 0 |
| 25 | 0 | 0 | 10 | 10 | 0 | 0 | 0 |
| 26 | 0 | 0 | 7 | 7 | 0 | 0 | 0 |
| 27 | 0 | 1 | 3 | 3 | 0 | .333 | 0 |
| 28 | 0 | 1 | 3 | 3 | 0 | .333 | 0 |
| 29 | 1 | 2 | 10 | 10 | .5 | .3 | 0.38 |
| 30 | 0 | 0 | 8 | 8 | 0 | 0 | 0 |
| 31 | 0 | 0 | 2 | 2 | 0 | 0 | 0 |
| 32 | 1 | 1 | 7 | 7 | 1 | .286 | 0.44 |
| 33 | 3 | 2 | 4 | 4 | 1.5 | 1.25 | 1.36 |
| 34 | 0 | 3 | 7 | 7 | 0 | 0 | 0 |
| 35 | 0 | 0 | 4 | 4 | 0 | 0 | 0 |
| 36 | 0 | 1 | 3 | 3 | 0 | .333 | 0 |
| 37 | 0 | 1 | 10 | 10 | 0 | .1 | 0 |
| 38 | 0 | 0 | 11 | 11 | 0 | 0 | 0 |

| | | | | | | | |
|----|---|---|----|----|----|-------|------|
| 39 | 1 | 0 | 14 | 14 | 0 | .07 | 0 |
| 40 | 0 | 1 | 12 | 12 | 0 | .083 | 0 |
| 41 | 2 | 0 | 4 | 4 | 0 | .5 | 0 |
| 42 | 0 | 0 | 6 | 6 | 0 | 0 | 0 |
| 43 | 0 | 0 | 7 | 7 | 0 | 0 | 0 |
| 44 | 0 | 0 | 4 | 4 | 0 | 0 | 0 |
| 45 | 1 | 5 | 13 | 13 | .2 | .43 | 0.27 |
| 46 | 0 | 1 | 16 | 16 | 0 | .0625 | 0 |
| 47 | 0 | 1 | 15 | 15 | 0 | .066 | 0 |
| 48 | 0 | 0 | 12 | 12 | 0 | 0 | 0 |
| 49 | 0 | 1 | 3 | 3 | 0 | .333 | 0 |
| 50 | 0 | 0 | 7 | 7 | 0 | 0 | 0 |

7. COMPARATIVE ANALYSIS OF PRECISION AND RECALL AFTER APPLYING PROPOSED ALGORITHMS AND FIND F-MEASURE

Here again it is analyzing the 1400 documents and on 50 queries of CRANTOP database. It is again analyzed on the proposed algorithm i.e. results of proposed algorithms which are found retrieved and also the relevant results in that how many are matched and how many are relevant results and on proposed algorithm how many are relevant results then we found the precision, recall and f-score now we were found that its values of precision, recall and f-score will be increased here after analyzed.

Table2. Demonstration of proposed algorithm through Precision, Recall and F-measure

| Q.n o | Results of propose algo. | Relevant result | Match ID's | Relevant result | Propose Algo Relevant Results | Precision | Recall | F-score |
|-------|--|---|------------|-----------------|-------------------------------|-----------|--------|---------|
| 1 | 12,13,2 4,29,47 ,51,141 ,172,18 4,431,4 97,519, 649,70 0,746,7 98,876, 879,88 0,948,9 76,119 7,1289, 1379 | 12,13,1 4,15,31 ,57,66, 95, 462, 497, 858, 879,88 0,184,2 9,51,10 2,378,8 59,185, 30,37,5 2,142,1 95,875, 56,876, 486 | 9 | 29 | 25 | 0.36 | 0 | 1.08 |
| 2 | 12,14,2 4,29,47 ,51,141 ,172,18 4,431,4 97,519, 649,70 0,746,7 98,876, 879,88 80,746, 859,94 | 12,14,1 5,51,10 2,202,3 80,390, 391, 658,87 7,948,1 84,52,3 80,746, 859,94 | 7 | 25 | 25 | 0.28 | 0 | 0.84 |

| | | | | | | | | | | | | | | | | | | | |
|----|--|---|---|---|----|--------------|--------------------------------------|----------|--|--|--|--|--|--|--|--|--|--|----------|
| | 0,948,9 76,119 7,1289, 1379 | 8,285,3 90,391, 442,64 3,858,4 97,856, 857,87 7,864,6 58,486, 858,64 3 | | | | | | | | | | | | | | | | | 15 79 |
| 3 | 5,6,90, 91,144, 168,18 1,399,4 85,518 | 5,6,181 ,144,48 5,90,91 ,119,39 9 | 7 | 9 | 10 | 0.7 | 0 . 7 7 7 7 8 | 2. 1 | | | | | | | | | | | |
| 4 | 236,40 1,1296, 1297 | 166,23 6,488 | 1 | 3 | 4 | 0.25 | 0 . 3 3 3 3 3 | 0. 75 | | | | | | | | | | | |
| 5 | 19,28,3 1,37,48 ,56,57, 68,85,9 5,122,1 60,272, 304,30 5,307,3 08,310, 360,40 1,495,4 97,525, 536,54 4,552,5 56,557, 570,57 2,976,1 076,12 72,137 9 | 401,48 8,1297, 522,12 96 | 1 | 5 | 34 | 0.02 9412 | 0 . 2 | | | | | | | | | | | | |
| 6 | 257,41 8,491,5 58,976 | 99,115, 257,25 8,491 | 2 | 5 | 5 | 0.4 | 0 . 4 | 1. 2 | | | | | | | | | | | |
| 7 | 48,56,5 7,58,35 4,469,4 92,717 | 19,20,5 6,57,58 ,492 | 4 | 6 | 8 | 0.5 | 0 . 6 6 6 6 6 7 | 1. 5 | | | | | | | | | | | |
| 8 | 48,56,5 7,58,35 4,469,4 92,717 | 20,999, 1005,4 8,122,5 8,196,3 54,360, 197,11 12,492 | 4 | 5 | 8 | 0.5 | 0 . 8 | 1. 5 | | | | | | | | | | | |
| 9 | 21,22,2 4,37,11 | 22,534, 21,550 | 4 | 4 | 19 | 0.21 0526 | 1 | 0. 63 | | | | | | | | | | | |
| 10 | 166,16 7,168,1 69,185, 236,40 5,427,4 88,493, 518 | 259,40 5,302,4 36,437, 438,90 8,1011, 493 | 2 | 9 | 11 | 0.18 1818 | 0 . 2 2 2 2 2 | | | | | | | | | | | | |
| 11 | 20,27,2 8 | 27,28,2 62,160, 20,263, 654,49 5 | 3 | 8 | 3 | | | | | | | | | | | | | | |
| 12 | 1,14,29 ,51,66, 141,14 2,172,1 85,272, 289,48 6,658,7 04,715, 749,85 9,877,8 99,999, 1197,1 272,12 89,133 3 | 86,194, 650,64 9,652,6 24 | 0 | 6 | 24 | 0 | 0 | | | | | | | | | | | | |
| 13 | 199,25 2,315,4 39,440, 467,46 8,469,4 96,503, 521,52 6,643,7 99,800, 879,90 3,919,1 290 | 64,265, 65,311, 496 | 1 | 5 | 19 | 0.05 2632 | 0 . 2 | | | | | | | | | | | | |
| 14 | 64,65,2 63,265, 517,60 9,903 | 64,65,4 96 | 2 | 3 | 7 | 0.28 5714 | 0 . 6 6 6 6 6 7 | | | | | | | | | | | | |
| 15 | 30,184, 195,46 2,463 | 463,46 2,497 | 2 | 3 | 5 | 0.4 | 0 . 6 6 6 6 6 7 | | | | | | | | | | | | |
| 16 | 2,3,4,3 | 266,10 | 1 | 4 | 8 | 0.12 | 0 | 0. | | | | | | | | | | | |

| | | | | | | | | |
|----|---|--|----------|-----------|----------|--------------|---|----|
| | 9,86,94 ,106,41 8 | 6,196,4 98 | | | | 5 | . | 37 |
| 17 | 8,15,94 ,106,19 9,242.2 59,279, 326,70 0,714,7 99,903, 1112 | 106,19 6,498 | 1 | 3 | 14 | 0.07 1429 | 0 | 0. |
| 18 | 52,56,5 7,58,85 ,121,12 2,173,1 96,197, 215,24 7,250,2 59,354, 360,40 9,469,4 92,498, 514,52 8,553,7 17,125 9 | 196,19 7,198,4 98 | 3 | 4 | 25 | 0.12 | 0 | 0. |
| 19 | 142,16 4,555,1 379 | 32,67,1 64,639, 715,71 6,719,1 379,71 7,499 | 2 | 10 | 4 | 0.5 | 0 | 1. |
| 20 | 87,88,2 68,269, 270,50 0 | 87,88,1 04,267, 268,26 9,270,4 07,408, 500 | 6 | 10 | 6 | 1 | 0 | 3 |
| 21 | 68,167, 185,23 6,283,3 02,405, 427,52 4,1355 | 271,16, 413,41 4,502 | 0 | 5 | 10 | 0 | 0 | 0 |
| 22 | 2,9,17, 23,132, 145,49 3 | 68,502 | 0 | 2 | 7 | 0 | 0 | 0 |
| 23 | 1,14,29 ,141,14 2,185,8 59,877, 899,11 12,119 7,1272, 1289,1 333 | 900,90 2,200,2 01,601, 899,90 3,593,1 99,594, 901,54 4,597,7 49,917, 919,13 33,634, 687,69 8,1290, 700,70 4,705,1 109,11 12,114 1,1197, | 5 | 33 | 14 | 0.35 7143 | 0 | 1. |

| | | | | | | | | |
|----|--|---|----------|-----------|----------|--------------|---|----|
| | | 1256,1 259,12 72,128 9,892 | | | | | | |
| 24 | 12,51,7 56 | 46,47,9 2,756 | 1 | 4 | 3 | 0.33 3333 | 0 | 0. |
| 25 | 7,31,48 ,52,95, 121,12 2,173,1 87,188, 200,20 1,212,2 13,214, 216,22 4,225,2 42,259, 272,27 6,278,2 79,391, 406,40 9,428,4 64,511, 513,51 9,880,1 197,12 59,127 2 | 213,21 2,214,2 15,216, 276,27 7,426,4 27,511 | 6 | 10 | 36 | 0.16 6667 | 0 | 0. |
| 26 | 3,4,8,1 1,16,79 ,94,104 ,145,20 7,265,2 71,306, 307,31 6,326,4 81,493, 526,52 7,528,5 65,570, 572,61 1,629,7 98,107 6,1109, 1185,1 226,13 55 | 145,61 1,376,4 06,565, 1076,5 11 | 4 | 7 | 32 | 0.12 5 | 0 | 0. |
| 27 | 1,30,19 5,199,2 01,225, 247,25 0,279,2 88,289, 433,49 7,520,6 43,652, 704,70 5,749,7 52,877, 901,91 9,1289, 1290,1 333 | 224,27 8,428,5 12 | 0 | 4 | 26 | 0 | 0 | 0 |
| 28 | 168,24 7,470,5 | 224,27 9,512 | 0 | 3 | 4 | 0 | 0 | 0 |

| | | | | | | | | |
|----|--|---|----------|-----------|----------|--------------|--------------------------------------|----------------------|
| | 93 | | | | | | | |
| 29 | 55,251, 464,46 5,466,6 01,612, 752,87 9,901,9 02,919, 1289 | 250,51 4,609,2 25,793, 464,46 5,612,4 66,513 | 4 | 10 | 13 | 0.30 7692 | 0 . 4 | 0. 92 30 77 |
| 30 | 123,25 0,289,4 64,465, 466,51 3,556,5 57,572, 752,90 2,919,1 197 | 225,46 4,514,4 66,609, 612,79 3,513 | 3 | 8 | 14 | 0.21 4286 | 0 . 3 7 5 | 0. 64 28 57 |
| 31 | 173,19 9,215,2 47,277, 433,48 7,512,5 97,652, 698,74 9,751,7 75,903, 919,10 76 | 776,75 1 | 1 | 2 | 17 | 0.05 8824 | 0 . 5 | 0. 17 64 71 |
| 32 | 4,250,3 02,306, 611,75 2,777,8 77,135 5 | 465,24 9,777,7 78,247, 250,75 2 | 3 | 7 | 9 | 0.33 3333 | 0 . 4 2 8 5 7 1 | 1 |
| 33 | 431,46 6,486,5 14,516, 544,55 3,594,6 09,612, 672,67 2,714,7 19,799, 800,87 6 | 252,43 1,141,5 16 | 2 | 4 | 17 | 0.11 7647 | 0 . 5 | 0. 35 29 41 |
| 34 | 431,46 6,486,5 14,516, 544,55 3,594,6 09,612, 672,71 4,719,7 99,800, 876 | 252,43 1,672,7 14,799, 800,51 6 | 6 | 7 | 16 | 0.37 5 | 0 . 8 5 7 1 4 3 | 1. 12 5 |
| 35 | 68,166, 167,16 9,185,4 05,488 | 166,16 7,132,5 17 | 2 | 4 | 7 | 0.28 5714 | 0 . 5 | 0. 85 71 43 |
| 36 | 13,21,2 9,51,66 ,119,14 2,269,2 | 169,16 8,518 | 0 | 3 | 22 | 0 | 0 | 0 |

| | | | | | | | | |
|----|---|--|----------|-----------|----|--------------|--------------------------------------|----------------------|
| | 70,283, 354,43 6,437,4 38,485, 493,52 4,550,5 54,623, 1185,1 226 | | | | | | | |
| 37 | 97,121, 173,18 8,242,2 77,409, 427,43 1,465,4 87,519 | 173,18 8,172,9 7,121,1 87,242, 409,48 7,519 | 8 | 10 | 12 | 0.66 6667 | 0 . 8 | 2 |
| 38 | 7,8,9,1 9,28,37 ,56,57, 68,85,7 9,122,1 60,187, 207,27 2,304,3 05,307, 308,31 0,314,3 15,360, 401,41 8,505,5 25,536, 544,55 6,557,5 58,570, 572,97 6,1076, 1272,1 379, | 24,283, 552,27 2,557,5 58,553, 554,55 5,556,5 36 | 5 | 11 | 39 | 0.12 8205 | 0 . 4 5 4 5 4 5 | 0. 38 46 15 |
| 39 | 2,3,4,7, 8,9,12, 16,17,2 1,22,23 ,37,79, 94,104, 145,20 7,265,2 71,306, 307,31 4,316,3 26,376, 406,48 1,493,5 26,527, 528,56 5,570,5 72,611, 629,79 8,899,1 076,11 09,118 5,1226, 1355 | 272,55 5,24,28 3,552,7 9,207,4 18,557, 554,55 6,505,1 257,53 6 | 2 | 14 | 44 | 0.04 5455 | 0 . 1 4 2 8 5 7 | 0. 13 63 64 |
| 40 | 7,8,9,1 9,28,37 ,56,57, 68,85,7 9,122,1 | 24,283, 552,27 2,85,97 6,557,5 58,553, | 7 | 13 | 39 | 0.17 9487 | 0 . 5 3 8 | 0. 53 84 62 |

| | | | | | | | | |
|----|--|---|----|----|----|--------------|---|---------------------------------------|
| | 60,187, 207,27 2,304,3 05,307, 308,31 0,314,3 15,360, 401,41 8,505,5 25,536, 544,55 6,557,5 58,570, 572,97 6,1076, 1272,1 379, | 554,55 5,556,5 36 | | | | | 4 6 2 | |
| 41 | 289,43 3 | 288,28 9,433,5 20 | 2 | 4 | 2 | 1 | 0 . 5 | 3 |
| 42 | 194,20 0,216,4 39,440, 467,46 8,469,4 96,503, 521,52 6,797,9 19 | 468,46 7,469,4 70,775, 521 | 4 | 6 | 14 | 0.28 5714 | 0 . 6 6 6 6 6 6 7 | 0 . 6 6 6 6 7 43 |
| 43 | 194,21 6,467,4 68,469, 496,50 3,521,5 26,797, 919 | 467,46 8,39,50 3,775,4 69,521 | 5 | 7 | 11 | 0.45 4545 | 0 . 7 1 4 2 8 6 | 1. 36 36 36 |
| 44 | 102,40 1,552,7 46,122 6,1296, 1297 | 302,43 6,437,5 24 | 0 | 4 | 7 | 0 | 0 | 0 |
| 45 | 2,14,30 5,308,3 10,525, 572,11 85 | 305,57 0,308,4 81,338, 1226,1 355,11 85,629, 663,79 8,572,5 25 | 5 | 13 | 8 | 0.62 5 | 0 . 3 8 4 6 1 5 | 1. 87 5 |
| 46 | 9,17,19 .28,37, 84,122, 123,16 0,304,3 05,338, 344,36 0,481,4 95,525, 544,57 2,623,1 076,11 85 | 305,33 8,344,4 81,84,1 23,118 5,623,5 70,798, 1226,6 29,663, 572,13 55,525 | 10 | 16 | 22 | 0.45 4545 | 0 . 6 2 5 | 1. 36 36 36 |
| 47 | 2,9,14, 22,87,2 | 304,30 6,307,3 | 9 | 15 | 17 | 0.52 9412 | 0 . | 1. 58 |

| | | | | | | | | |
|----|---|--|---|----|----|--------------|--------------------------------------|----------------------|
| | 07,304, 305,30 6,308,3 09,310, 464,52 5,572,6 11,663 | 05,629, 663,30 8,309,3 10,570, 798,11 85,135 5,572,5 25 | | | | | 6 82 | 35 |
| 48 | 2,8,9,1 97,200, 304,30 7,308,3 10,315, 439,44 0,464,5 11,513, 514,52 1,526,5 44,572, 593,60 1,634,7 97,799, 800,87 9,880,9 00,901 | 439,31 1,316,4 40,187, 314,31 5,797,7 98,794, 265,52 6 | 5 | 12 | 30 | 0.16 6667 | 0 . 4 1 6 6 6 7 | 0. 5 |
| 49 | 23,320, 478,52 7 | 320,47 8,527 | 3 | 3 | 4 | 0.75 | 1 | 2. 25 |
| 50 | 106,18 8,196,1 97,279, 326,46 8,498,5 28,550, 714,11 12,125 9 | 326,62 9,21,22 ,94,306 ,528 | 2 | 7 | 13 | 0.15 3846 | 0 . 2 8 5 7 1 4 | 0. 46 15 38 |

8. CONCLUSION

Comparative Analysis of very simple intersection algorithms have been done on queries, which have been picked from the standard crantop database. More advanced algorithms are in use, but the simple algorithms have been picked to emphasize the need of a new algorithm by showing the computed values of Precision, Recall and F-Score. These values are very low, depicting the need to develop a novel idea for searching documents. After applying the advanced proposed algorithm the values of Precision, Recall and F-Score will become high as compared with previous algorithms for the information retrieval. The significance of new algorithm is that its database can take root words, Break Query, Phrase Queries, Remove Stopwords, Convert Remaining Into Root Words, Classify Into Specific And General Words, Synonyms, Intersection With Specific Words by doing so it can increase the value of Precision, Recall and F-Score if we are using sparse matrix structure to do that.

9. REFERENCES

- [1] Alexander B. et al., 2009 "A decision theoretic approach to combining information filtering", Online Information Review, Vol. 33 Iss: 5, pp.920 – 942
- [2] E. Michael Keen, 1991 "The use of term position devices in ranked output experiments", journal of documentation, vol. 47 iss: 1, pp.1 – 22

- [3] François X. D. et al.,2014 "applying and theorizing institutional frameworks in is research: a systematic analysis from 1999 to 2009", *information technology & people*, vol. 27 iss: 3, pp.280 – 317
- [4] Judit b. et al., 2006 "methods for evaluating dynamic changes in search engine rankings: a case study", *journal of documentation*, vol. 62 iss: 6, pp.708 – 729
- [5] Kevin C. O., 1996 "world wide web-based information storage and retrieval", *online and cd-rom review*, vol. 20 iss: 1, pp.11 – 20
- [6] Peter I., 1996 "cognitive perspectives of information retrieval interaction: elements of a cognitive ir theory", *journal of documentation*, vol. 52 iss: 1, pp.3 – 50
- [7] Soren B., 1996 "Cybersemiotics: a new interdisciplinary development applied to the problems of knowledge organisation and document retrieval in information science", *journal of documentation*, vol. 52 iss: 3, pp.296 – 344
- [8] S.M. Zayed Ahmed et al., 2009 "A review of research on human-computer interfaces for online information retrieval systems", *The Electronic Library*, Vol. 27 Iss: 1, pp.96 – 116
- [9] Stella D. T. et al. ,2014 "Using cues to forage for information on the Web", *Journal of Systems and Information Technology*, Vol. 16 Iss: 4, pp.296 – 312
- [10] Valery J. F. et al., 1997, *Control and feedback in ir systems*, in valery j. Frants, jacob shapiro, vladimir g. Voiskunskii (ed.) *Automated information retrieval: theory and methods (library and information science, volume 97a)* emerald group publishing limited, pp.222 – 259
- [11] ConvertedtoXMLfrom:ftp://ftp.cs.cornell.edu/pub/smart/cran/ (for 1400 documents)
- [12] ConvertedtoXMLfrom:ftp://ftp.cs.cornell.edu/pub/smart/cran/(for 50 queries)
- [13] EllisHorowitz, Sartaj Sahni, Dinesh Mehta, *Fundamentals of data structures in C++,2nd Edition* ,2008, Universities Press(India) Private Limited.
- [14] Christopher D. Manning(Stanford University) ,Prabhakar Raghvan(Yahoo! Research), Hinrich Schutze(University of Stuttgart) , *Introduction to Information Retrieval*,2009,Cambridge University Press.