



Comparison of Term Weighting Models for Retrieval of Information

Prachi Garhewal

Research Scholar

Shri Ram Group of Institutions (SRGI)
Jabalpur (M.P.), [INDIA]

Email: gahrhewalprachi@gmail.com

Anupam Choudhary

Lecturer,

Kalaniketan Polytechnic College
Jabalpur (M.P.), [INDIA]

Email: chowdharyanupam7@yahoo.com

Santosh Kumar

Assistant Professor

Shri Ram Group of Institutions (SRGI)
Jabalpur (M.P.), [INDIA]

Email: sant303@rediffmail.com

Abstract—Information Retrieval is finding documents of unstructured nature which should satisfy user's information needs. There exist various models for weighting terms of corpus documents and query terms. This work is carried out to analyze and evaluate the retrieval effectiveness of various IR models while using the new data set of FIRE 2011. The experiments were performed with *tf-idf* and its variants along with probabilistic models. For all experiments and evaluation the open search engine, Terrier 3.5 was used. Our result shows that *tf-idf* model gives the highest precision values with the news corpus dataset.

Keywords:—TF-IDF, BB2, Fire Dataset, Retrieval Effectiveness, Precision, Recall, Information Retrieval, IR Models, Weighting Schemes

query of user, according to that gives the ranks to the all documents and bring the top relevant documents from data set. For this paper we use static data set to evaluate the result of different models. Before assigning ranks to the documents, information retrieval system goes through some preprocessing steps that will be discussed in the next section of the paper.

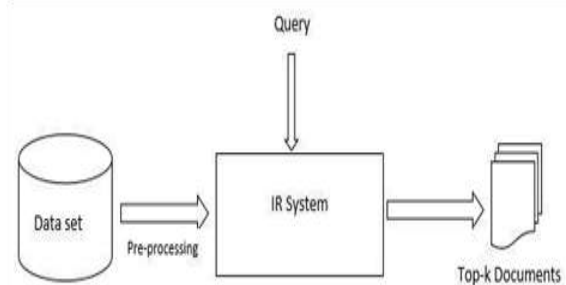


Figure 1: Information Retrieval

1. INTRODUCTION

Information retrieval is finding documents of unstructured nature which should satisfy user information need [1]. The objective of information retrieval system is to retrieve the documents from the collection of documents that fulfill the user need which is express in terms of query. Information retrieval system uses many models to understand the

2. PRE-PROCESSING

Every document must go through some preprocessing steps, to make it simple for search engines to take them as a document and run them on their algorithms[2]. Steps like tokenization, stemming, normalization, stop word removing etc. are applied on each document.

2.1 Tokenization

Tokenization is define as the process of breaking text documents into words, phrases, numbers, symbols etc. all these are called tokens. Tokenization faces issues like language recognition, as this process is language dependent. User data, Meta data, machine learning methods are useful to determine information about document language.

2.2 Stemming

Stemming is the process in which reducing the words by removing letters to their root word. A root word can represent many words that might have different meaning and root word may be doesn't make any meaning some times. Avoiding the many original words and using stem words helps to reduce the size of dictionary, that contain all words of document collection. In other sense stemming help user by providing the choices of related options of user query, just after typing the few letter in query box.

2.3 Stop-words

A document contains hundreds or thousands words, for the user perspective every word is not equally important. Generally this word appears many times in the document and doesn't contribute any information for the document makes it informative. As we think practically, people don't search words like 'the', 'a', 'of' and many other words these words called stop words. Stop-words list are can different for different document collection based on the purpose.

2.4 Inverted Index

Every documents collection contains many documents and a document has many different words. Now question is how we make search very fast and how to know which document contain query terms. The answer is inverted index, it is kind of link list in which every word from whole collection connected with the nodes that represent documents numbers in which that particular term appears.

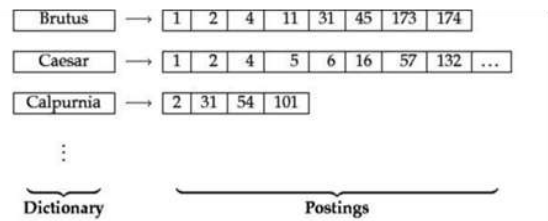


Figure 2: Inverted Index [1]

Generally we use stemmed words in the index, like compute for computer, computation, computations and many other similar words as well doesn't include stop word in index, reason is just to avoid space problem for system. In above figure, collection of all searchable words called dictionary and link of documents IDs called posting.

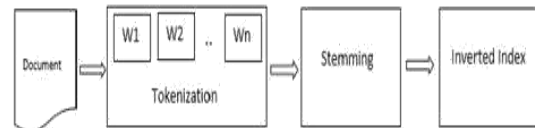


Figure 3: Pre-Processing

3. REVIEW OF IR MODELS

Information Models define the way to represent the document text and the query as well. Other objective of IR models is to compare the document and query, assign ranks to documents. We are using many version of different model such as Probabilistic models, TF-IDF weighting model and Divergence from randomness.

3.1 TF-IDF Model

Generally Term frequency(TF) define as how many times a term appears in a document and document frequency known as in how many documents a term appears.

3.2 BB2 Model

BB2 Model BB2 model is Bernoulli-Einstein model with Bernoulli after-effect and normalization is vector space model that use key frequency, document length, norm to score assigned document.

4. EXPERIMENTAL EVALUATION

Evaluation is always an important part of any research of any area in all around the world, same it useful in context of information retrieval. Evaluation in the simple mean is how effective a system performs and produces valuable result with accuracy.

4.1 Evaluation Measures

In Information Retrieval, to measure the effectiveness of the system our requirement is a data set, a set of queries and some function to judge relevance factor between document and queries. Simple IR system just fetches the best relevant documents that are related to the query and assign ranks to them. Now effectiveness depends on measurements used for evaluation, better measurements give better ranked list of documents. In this paper we use very common measurements such as precision and recall that are discussed in next section.

4.1.1 Precision

In simple words, precision can be defined as the ratio of number of relevant retrieved documents to the number of retrieved documents [8].

Precision generally mention in form of percentage and as the number of retrieved documents increase, the precision of system will decrease.

4.1.2 Recall

Recall is another measure for information Retrieval model, which can be described as a ratio of number of relevant retrieved documents to the number of relevant documents [8].

Recall and precision is inter depended, recall will increase when relevant retrieved documents increase. Recall and precision are inversely related.

4.2 Description of Data Set

The experiments has been carried out on the data set of FIRE 2011[http://www.isical.ac.in/~fire/] for English. The data set contains various documents from English news domain-The Telegraph. These news articles are extracted from 2001 to 2010 and contain 303,292 documents. We just took a sample form this data set for our experiment. The task of corpus creation was carried out to support experiments for research purpose in information retrieval domain.

4.2.1 FIRE and Document Format

FIRE is stands for form of Information Retrieval and Evaluation. It's an India based organization for research on information retrieval. FIRE works on languages of South Asian contraries.

Document format that used in FIRE collection follow the standard representation of TREC collection. Documents contain tags like DOC, DOCNO and TEXT. DOCNo is unique number for every document in the data set. Text field contains the actual news article in plain text. The example of a text file is shown below.

```
<DOC>
<DOCNO>doc_03/0003</
DOCNO> <TEXT>

Americans used more health services and
spent more on prescription drugs in 2013, re-
versing a recent trend, though greater use of
cheaper generic drugs helped control spending,
according to a report issued on Tuesday by a
leading healthcare information company.

</TEXT>

</DOC>
```

Figure 4. Document Format

4.2.2 Topic File

Topics file contain some pre-fixed queries for the data set, these queries almost

cover every document within the data set. According to our sampled data of FIRE data collection we take 9 queries. Example of our topic file is shown in figure 5. The topic file format contain tags such as top, num and title. Title is the query and number is assign to every topic.

```
<topics>
<top>
<num>1</num>
<title>infini thoughts
</title> </top>

<top>
<num>2</num>
<title>umesh publishers </title> </top>
```

Figure 5: Topic File

4.2.3 Qrels File

Qrels file format describes the presence and absence of the every query terms in the document. Format of qrels shown in figure 6 and description of format is like this, a first column show the query ID that is according to the topic file, second place show iteration, third place the document ID that is mention in document format and last column shows the presence and absence of that query in document by 0 or 1.

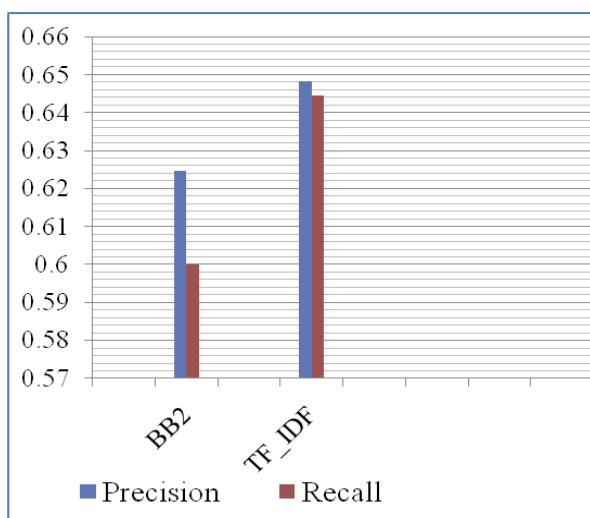


Figure 6: View of Qrel file

5. RESULT AND ANALYSIS

We performed our experiments in Terrier 3.5. It has all the necessary codes to support experiments for FIRE dataset. We make some changes in terrier. properties file. There is many Information model already supported by the terrier-3.5. we are showing the result of BB2 and TF-IDF. We use two measures Precision and Recall for comparing the models. In figure 5.1 example of eval file that generated for every model by terrier-3.5 and it shows information about retrieved and relevant documents. We applied various models in our dataset and compare the results. Table illustrates the result of comparisons. TF_IDF gives the Precision value of 0.6481 and it is greater than BB2 Model.

6. CONCLUSION

This work has been carried out to implement various Information Retrieval Models with the FIRE dataset which contains corpus of newspapers. We implemented the tf-idf model and BB2 Model and compare their results. Based on our results we conclude that tf-idf produces the better results than BB2 Model. The results were evaluated and successfully compared with Terrier.

REFERENCES:

- [1] An Introduction to Information Retrieval Christopher D. Manning Prabhakar Raghavan Hinrich Schütze.
- [2] Sager, Juan C. *A practical course in terminology processing*. John Benjamins Publishing, 1990.
- [3] Baeza-Yates, Ricardo, and Berthier Ribeiro-Neto. *Modern information retrieval*. Vol. 463. New York: ACM press, 1999.
- [4] Frakes, William B. "Stemming Algorithms." (1992): 131-160.
- [5] Patel, B. N., Prajapati, S. G., & Lakhtaria, K. I. (2012). Efficient Classification of Data Using Decision

- Tree. *Bonfring International Journal of Data Mining*, 2(1), 06-12.
- [6] Xia, Tian, and Yanmei Chai. "An Improvement to TF-IDF: Term Distribution based Term Weight Algorithm." *Journal of Software* (1796217X) 6.3 (2011).
- [7] Alvarez, Sergio A. "An exact analytical relation among recall, precision, and classification accuracy in information retrieval." Boston College, Boston, Technical Report BCCS-02-01 (2002): 1-22.
- [8] Akhilesh Sharma, Kamaljit Lakhtaria, Santosh Vishwakarma, "Data Mining Based Predictions For Employees Skill Enhancement Using Pro-Skill-Improvement Program & Performance Using Classifier Scheme Algorithm", *International Journal of Advanced Research in Computer Science*, ISSN No. 0976-5697, Vol. 4, No. 3, March 2013, Page No. 102 – 107.
- [9] Robertson, Stephen. "Understanding inverse document frequency: on theoretical arguments for IDF." *Journal of documentation* 60.5 (2004): 503
- [10] Santosh K. Vishwakarma, Kamaljit I Lakhtaria, Divya Bhatnagar, Akhilesh Sharma (2014). "An efficient approach for inverted index pruning based on document relevance" *Conference Proceeding of Fourth International Conference on Communication Systems and Network Technologies*, Page No. 487 -DOI 10.1109/CSNT.2014.103
- [11] Lakhtaria, Kamaljit I., Bhaskar N. Patel. "Implementing R-Tree Index Optimization in Core Banking system." *International Journal of Research in Management, Economics & Commerce*, 2(3) (2012), 42-48
- [12] Saracevic, Tefko. "Evaluation of evaluation in information retrieval." *Proceedings of the 18th annual international ACM SIGIR conference on Research and development in information retrieval*. ACM, 1995.
- [13] Amati, Giambattista. *Probability models for information retrieval based on divergence from randomness*. Diss. University of Glasgow, 2003.
- [14] Robertson, Stephen, and Hugo Zaragoza. *The probabilistic relevance framework: BM25 and beyond*. Now Publishers Inc, 2009.
- [15] Lakhtaria, Kamaljit I. *Technological Advancements and Applications in Mobile Ad-hoc Networks: Research Trends*. Information Science Reference, 2012.
- [16] Lakhtaria, K. I., Patel, P., & Gandhi, A. (2010). *Enhancing Curriculum Acceptance among Students with E-learning 2.0*. arXiv preprint arXiv:1004.2560.
- [17] www.terrier.org
- [18] Sharma, Akhilesh K., Kamaljit I. Lakhtaria, Avinash Panwar, and Santosh K. Vishwakarma. "An efficient approach using LPFT for the karaoke formation of musical song." In *Advance Computing Conference (IACC), 2014 IEEE International*, pp. 601 - 605. IEEE, 2014.