



## Information Retrieval by Local Value schemes: A Review

Radha Joshi<sup>1</sup> and Basant B Joshi<sup>2</sup>

<sup>1</sup>Dept. of Applied Science, Amrapali Institute of Technology, Haldwani (Uttarakhand), India

<sup>2</sup>Associate Analyst, Content Engineer, Global Logic Technologies, Gurgaon (India)

**ABSTRACT:** Information Retrieval system's performance depends its value scheme. Value can be seen in one aspect that is local. For each type of value scheme single term is considered. Sizing and term dependency is natural in documents. In the present paper an attempt has been made to study the dependency of term and value to each other. Term dependency has been use to define local value scheme.

**Keywords:** Dependency ,sizing ,information retrieval

### I. INTRODUCTION

Best example of seeing the information Retrieval is internet search engine. When we try to search something on search engine lot of information retrieval indexed open manually. Search engine are working on the formula of Boolean system, vector space model, binary model and local scheme. Here we focus on valuescheme. Localvalue are function of frequency of database index. Anyone can easily find out the information Retrieval technique form database.in Manning & Schatze and sander its described in their book[7].The process of retrieve data from backend to search engine can be seen in Chishlm & Kolda [1].Different type of data retrieving technique based on their size are discussed on lonet [8]. Here we will discuss some technique of information retrieval-

**Binary technique-**The best technique to retrieve the information is binary technique this technique is expressed below-

$$v_{ij} = \begin{cases} 1 & \text{if } f_{ij} > 0 \\ 0 & \text{if } f_{ij} = 0 \end{cases} \quad (1)$$

Where  $V_{ij}$  is term frequency of term  $i$  in any document  $j$ .

Frequency value scheme-This shows the frequency of term

$$v_{ij} = f_{ij} \quad (2)$$

Frequency's Logarithms-Frequency term is not good for any information retrieval, more frequency means more dependency of data, here we need to reduce the dependency of data into each other. Forexample, any term  $i$  comes 30 times in any information and any term  $j$  comes 2-3 times then according to frequency method term  $i$  is more important than but in real life practice it is not necessary. But when we use the logarithm of the frequency it becomes more easy to retrieve the information regarding that frequency. This can be defined as-

$$V_{ij} = \begin{cases} \log(1 + f_{ij}) & \text{if } f_{ij} > 0 \\ 0 & \text{if } f_{ij} = 0 \end{cases} \quad (3)$$

Normalization is one of the important term for calculate the repeating frequency on a fixed time interval [0 1]. With the help of normalization we can divide the whole table into one or more than one table for the better retrieval we can also generate the sub table or we can say pivot table. As on maximum  $f_{ij}$  where  $1 \leq i \leq n$  is normalized factor for frequency repeating , if there are  $n$  total number of document.

Repeating value schemes always depend on the previous one entry. Repeating schemes based on the assumption of independence of term to term informed documents. Kim Hee-Soo *et al* [2] computed repeating dependence. In place of single search its better to describe in paired documents.

Term frequency and documents length combination is used to infer a saliency of a term in a documents, sometimes it's very good to choose only term frequency because term frequency are very small and it's interesting to emphasize the frequencies in a documents.

In the present review paper, we have discussed some local repeating scheme for pair which performs much faster than in single word search.

### II. PROPOSED REPEATING SCHEMES FOR WORD PAIRS FOR FASTER PERFORMANCE-

We have purposed some steps for retrieve the information, which is based on individual value of words and the world pair which is very close to each other. That steps are given below-

1. Select a document for repeating value scheme.
2. Remove all words from given article, determine propositions and helping etc.
3. Change every word to its lexical form, for example demonization and demonetizations with demonetize.

4. Count the Repeating frequency of each function. Any word  $v_1$  count frequency will-  
 $f(v_1)$ =number of statement where  $v_1$  appears.

5. Select the most repeating word.

6. Now make the combinations of most frequent words with word pair and count the frequency of word pair.

Combination of each frequency count as :

$f(v_1, v_2)$  = number of statement in which both  $v_1, v_2$  appear.

7. Now calculate the probability of  $(v_1, v_2)$  word pairs in relevant documents and this relationship can be denoted as  $P_r(v_1, v_2)$  and defined by

$$P_r(v_1, v_2) = \frac{f(v_1, v_2)}{\text{Min}[f(v_1), f(v_2)]} \quad (4)$$

8. Calculate the repeating pair of  $(v_1, v_2)$

$$V(v_1, v_2) = [\log(1 + f(v_1)) + \log(1 + f(v_2))] * P_r(v_1)(v_2) \quad (5)$$

Proposed repeating scheme is a combinations of logarithm of frequency of each word multiplied by repeating of word pair in the form of probability.

### III. EXPERIMENTAL OUTPUT

For the experiment we choose some common or famous trend text, then we will divide their output. Some famous text trend are - 'modi india pm', second 'modi the public servant' and third - 'real reason of demonization' repeating frequency of word pair have been calculated

#### Document 1

$v_1$	$v_2$	$f(v_1, v_2)$	$W(v_1, v_2)$
India	PM	8	2.04
Modi	demonetization	6	1.21
Modi	CM	14	0.81
CM	Yogi	2	0.76
Modi	Gujrat	5	0.2
UP	Yogi	1	0.31

#### Document ii

$v_1$	$v_2$	$f(v_1, v_2)$	$W(v_1, v_2)$
Modi	India	10	2.66
Modi	PM	7	2.28
Modi	CM	7	1.52
Modi	Gujrat	4	1.31
Modi	demonization	2	1.31
Modi	Skill India	2	0.81
Modi	BJP	2	1.50

#### Documents iii

$v_1$	$v_2$	$f(v_1, v_2)$	$W(v_1, v_2)$
India	PM	3	0.96
Demonization	Modi	1	0.52
Demonetization	Russia	1	0.51
Demonetization	Economy effect	1	0.34
Demonetization	Benefit	1	0.21
Demonetization	RBI	1	0.20
Modi	Japan	1	0.19

Documents First is related about Modi and his political culture and Second document about Modi and his major task and current scenario, last documents completely based on demonetization which was taken by PM Modi. In first document some highest frequency word are Modi -53, yogi-47, CM-33, demonitization-21.

In second documents frequent word are Modi-38, Demotization-27, PM-18, skill India-12, and in the last document most frequent word are demonetization -42, Modi-31, Russia-15, RBI-10. From first table it can be seen that CM and demonetization is low frequency word yet they make an important combination with other word. In 3<sup>rd</sup> document the word pair PM and demonetization has highest value which itself proves the content of the documents. Value of word pair show the dependency of word in the pair of each other.

### IV. CONCLUSION

In the present review paper, it has been shown that word pair can also describe a document. Though single term can describe the information of documents, but word pair gives clear idea about the content. The first two documents are related to Modi and frequency of two words are almost similar. In this case combination of text 'Modi' play important role with another word. Table 1<sup>st</sup> show that similar word pairs have higher value then other word pairs and have the capacity to describe the whole documents in a single manner. The review value can also be extended to a combination of more than three word or a complete phrase.

### REFERENCES

- [1]. Chisholm, E. and Kolda T. G. 1999. New term weighting formulas for the vector space method in information retrieval. Technical report ORNL-TM-13756, Oak Ridge national laboratory, Oak ridge, TN
- [2]. Dominich S. 2008. The Modern Algebra of Information Retrieval (Information Retrieval Series). Springer-verlag, New York.

- [3]. G. Salton and C. Buckley. Term weighting approaches in automatic text retrieval. *Information Processing and Management*, **24**(5): 513-523, 1988.
- [4]. W. B. Croft and D. J. Harper. Using probabilistic models of document retrieval without relevance information. *J. Documentation*, **35**(4): 285-295, 1979. Cited in [6].
- [5]. S. Butcher, C. L. A. Clarke, and I. Soboroff. The trec " 2006 terabyte track. In TREC, 2006.
- [6]. C. L. A. Clarke, N. Craswell, I. Soboroff, and E. M. Voorhees. Overview of the trec 2011 web track. In TREC, 2011
- [7]. Manning C. D. and Schutze H. 2002. Foundations of Statistical Natural Language Processing. MIT press, Cambridge.
- [8]. Lan M., Sung S. Y., Low H. B. and Tan C. L. 2005. A comparative study on term weighting schemes for text categorization. In *Proceedings of the International Joint Conference on Neural Networks (IJCNN)*, Vol. **1**, 546-551.
- [9]. Term weighting schemes and formula in information retrieval by Erica Chisholm and Tamara G. Kolda.