

# Improving Document Vectors Representation using Semantic Links and Attributes

**Chirag Shah**

Dept. of CSE

IIT Bombay, Powai

Mumbai 400076

Maharashtra, India.

Email: *chirag@iitb.ac.in*

**Pushpak Bhattacharyya**

Dept. of CSE

IIT Bombay, Powai

Mumbai 400076

Maharashtra, India.

Email: *pb@cse.iitb.ac.in*

## Abstract

Document representation is a crucial step in any Information Retrieval (IR) system. Since most of the traditional methods do not consider much of semantic or syntactic information, the representation becomes insufficiently informative for an IR task. We describe a novel approach to incorporating Natural Language Processing (NLP) in document representation for addressing this problem. Use is made of additional information about the sentences, *viz.*, (i) syntactic links among the words found by the *Link Parser* and (ii) heuristically determined semantic attributes of the words. After mapping this information to the document level using Self-Organizing Map (SOM), we use it for embellishing the document vectors constructed by the TFIDF method. The efficacy of the proposed method is established by showing that the document vectors (i) have higher *mutual information content* and (ii) achieve better class separation.

## 1 Introduction

The use of syntax and semantics for information retrieval (IR) is well studied in the literature. While many researchers accept the fact that in principle, such additional information should help in improving IR [Monz, 2000], there are divergent views on whether and how to use this information. For instance, Bowen Hui [Hui, 1998] examined some limitations of traditional IR systems and laid out the motivation for applying NLP techniques to IR, focusing at the same time on only morphological processing and stressing that "*not all natural language phenomena apply to IR*". Now the question is what phenomena are important.

In this paper we investigate the effect of some of such phenomena like link information and word attributes on IR. It is obvious that the *better* the documents are represented in terms of vectors, the better will be the performance of an IR system built using them. We, therefore, focus on representing the documents in a *better* way using with some NLP techniques in this paper. We argue that

*incorrectness* and *insufficiency* of the information are two major drawbacks of traditional document representation schemes and focus on addressing the problem of *insufficiency* of the information. The rest of the paper is mainly organized in four parts. Part I sets the motivation behind using NLP for IR in general and link and attributes information in particular. Part II describes our method for collecting additional information about the documents using Link Grammar [Sleator and Temperley, 1993] and other heuristics. This part demonstrates how relations and attributes information can help in supplying some additional and useful information about the documents. In order to measure the *goodness* of generated document vectors, we follow the intuition and compelling experimental evidences provided by Rong Jin *et al.* [Jin *et al.*, 2001] that *the more informative the document vectors are, the better will be the performance of IR using these vectors*. We, therefore, find the *informativeness* of the document vectors, which is explained in part III of the paper. In part IV we provide additional support for our method by finding interclass distance for traditional method and our proposed method. In section 7 we conclude that careful use of NLP can indeed improve performance of an IR system.

## I Motivation behind using NLP for IR

In this part we analyze the shortcomings of traditional methods of document representation and set the motivation for using NLP for IR in general, and document representation in particular. Here we are considering vector space model [Salton *et al.*, 1975], which is the most accepted and widely used method for document representation.

## 2 Shortcomings of the Traditional Methods

It is expected that the representation of documents should reflect the knowledge meant to be conveyed by the documents. The traditional methods for representation of documents like Term Frequency (TF) [Salton, 1989], Term Frequency with Inverse Document Frequency (TFIDF) [Joachims, 1997], Weighted IDF (WIDF) [Tokunaga and Iwayama, 1994] *etc.* do not consider the senses of the words or their mutual semantic relations. This causes problems. In general, the vector space model of document representation [Salton *et al.*, 1975] using *bag of words*, suffers from two problems:

### 1. Incorrectness of information

*Synonymy* (more than one word having the same sense), and *polysemy* (single word having more than one sense) affect recall and precision respectively [Deerwester *et al.*, 1990]. There are many studies for solving these problems using Word-Sense Disambiguation (WSD) [Mihalcea and Moldovan, 1999; Agirre and Rigau, 1996; Mihalcea and Moldovan, 1998]. However, [Sanderson, 1994] reports that disambiguation accuracy of at least 90% is required to avoid the degradation of effectiveness of retrieval.

### 2. Insufficiency of information

Many times merely considering the words of a document may not be enough for representation as a document is not a bag of words. For example, consider two documents using the same set of words, but one talks in a positive sense, while the other talks in a negative sense. Since the traditional methods of document representation do not have the mechanism to capture the tone or the structure of the document, both of these documents will get almost the same vector

representation. This results in coarse clustering or poor precision when such vectors are used in an IR system. It is obvious that a document is represented more appropriately if syntactic and semantic information is included. The present paper is concerned with this issue of insufficiency of information.

## II Using Relations and Attributes for Extending Document Vectors

Consider these following sentences: *John plays football. John did not play football. John likes to play football. John cannot play football. He asked John to play football.*

Traditional frequency based methods represent all these sentences<sup>1</sup> in the same way. However, there are differences among these sentences with respect to the tense, intention or ability of the agent, positive or negative connotation, semantic roles of the words and so on. These details demand a deeper analysis of the text. To facilitate this, we use Link Grammar [Sleator and Temperley, 1993] as one of the tools, which provides the link information among the words of a sentence. For example, for the sentence *John plays football*, the Link parser outputs

```

+--Ss+-----0s---+
|      |           |
John plays.v football.n

```

where *John* is the subject of *play* and *football* the object. Additionally, we make use of heuristics that capture semantic attributes. Table 1 show the list of attributes along with their explanations.

| Attribute | Meaning   |
|-----------|---|
| not       | Negative sense  |
| present   | Present tense   |
| past      | Past tense  |
| future    | Future tense  |
| def       | Definite  |
| indef     | Indefinite  |
| contrast  | Shows contrast in statements                                    |
| ability   | Demonstrates ability of some act                                |
| should    | To do something as a matter of course                           |
| may       | Possibility that something is true or happens                   |
| just      | Expresses an event or a state that has just begun or ended      |
| yet       | Expresses an event or a state that has not yet begun or started |
| progress  | An event is in progress   |
| request   | Request for something   |

Table 1: Semantic attributes and their meanings

<sup>1</sup>All words other than *John*, *play*, and *football* are eliminated as stop words [Yang and Wilbur, 1996].

### 3 Forming Sentence Category Map (SCM)

Since extracting links information and heuristic analysis is done at the sentence level, we need to extend it to the document level in order to obtain document representation. We make use of Self-Organizing Map (SOM) for this purpose. SOM falls under a special class of neural networks based on *competitive learning* [Haykin, 1995]. In such neural networks, output neurons of the network compete among themselves to be activated or fired, with the result that only *one* output neuron, or one neuron per group, is on at any one time. An output neuron that wins the competition is called a *winner takes all neuron* or simply a *winning neuron* [Yegnanarayana, 1999]. Various characteristics like approximation of the input space, topological ordering, density matching, and ability to select the best features, make SOM ideal for performing tasks like pattern organization [Haykin, 1995].

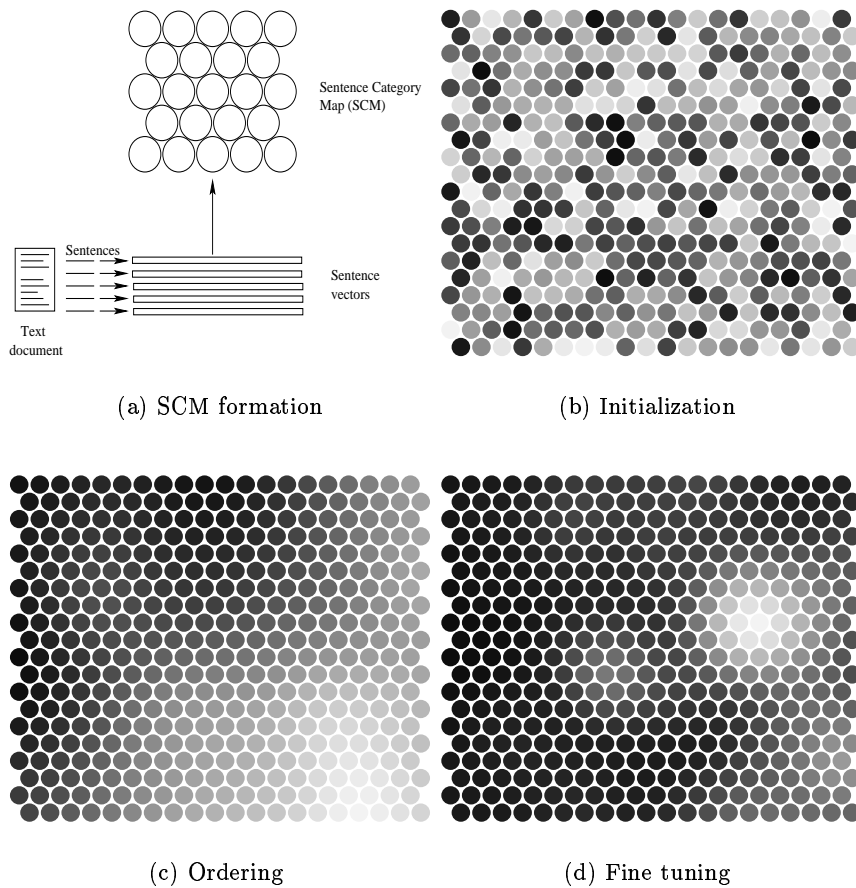


Figure 1: Forming Sentence Category Map (SCM)

In order to use SOM for organizing sentences, we use the modified version of Kohonen's algorithm

[Kohonen, 1995] and WEBSOM architecture [Honkela *et al.*, 1996; Kohonen *et al.*, 2000] as given below.

1. *Input patterns generation*: There are 107 Link Grammar relations and 14 attributes that we generate. Therefore, form a vector of size 121 for each sentence assigning. Each component of this vector will denote the count for the corresponding link or attribute in that sentence.
2. *Initialization*. Select the size of the SOM. We chose 20x20 (400 neurons). With each neuron  $j$ , there will be a weight vector  $w_j$  associated. Choose random values for the initial weight vectors  $w_j(0)$ . The only restriction here is that the  $w_j(0)$  be different for  $j = 1, 2, \dots, l$ , where  $l$  is the number of neurons in the lattice. It may be desirable to keep the magnitude of the weights small.
3. *Sampling*. Draw a sample  $X$  from the input space of  $N$  vectors.
4. *Similarity Matching*. Find the best-matching (winning) neuron  $i(X)$  at time step  $n$  by using the Euclidean minimum distance criterion

$$i(X) = \arg \left( \min_j \|X(n) - w_j(n)\| \right), j = 1, 2, \dots, l \quad (1)$$

5. *Updating*. Adjust the weight vectors of all neurons by using the update formula

$$w_j(n+1) = w_j(n) + \eta(n)h_{j,i(X)}(n)(X(n) - w_j(n)) \quad (2)$$

where  $\eta(n)$  is the learning-rate parameter, and  $h_{j,i(X)}(n)$  is the neighborhood function (we used *Gaussian*) centered around the winning neuron  $i(X)$ ; both  $\eta(n)$  and  $h_{j,i(X)}(n)$  are varied dynamically during learning for the best results as shown in [Haykin, 1995].

6. *Continuation*. Continue with steps 3, 4, and 5 until no noticeable changes in the feature map are observed.

This scheme is shown in Figure 1(a). We chose to train a SOM of size  $20 \times 20$ . The results of various stages of the algorithm are shown in Figures 1(b), 1(c), and 1(d). These figures represent the effect of input patterns on the SOM. Since initially the weight vectors corresponding to neurons are initialized to some random values, the input patterns get mapped at any position on the map. As the ordering of weight vectors according to the input patterns takes place, we can see a gradual organization in the map. More details about this procedure can be found in [Kohonen, 1995; Haykin, 1995].

After the map is constructed, we can observe a kind of organization of sentences according to their semantic representation given by the relations and the attributes. Therefore, we call it *Sentence Category Map (SCM)*.

## 4 TFIDF Vectors Enhancement with SCM

Once the sentences are organized using SOM, we can find the representation for the documents using the following algorithm.

1. For every document do the following. Take every sentence's vector and input it to the trained SCM. Find the winning neuron. This constitutes one *hit* on that neuron.

2. Collect hit information from all the neurons and construct a vector using it. In our case this vector is of size 400.
3. Append this vector to the normal TFIDF vector for the document.

This scheme is shown in Figure 2. The resulting vectors incorporate *additional* and *useful* information about the document. They are expected to give *better semantic representation* of the document. This *goodness* is justified in the next two parts of the paper.

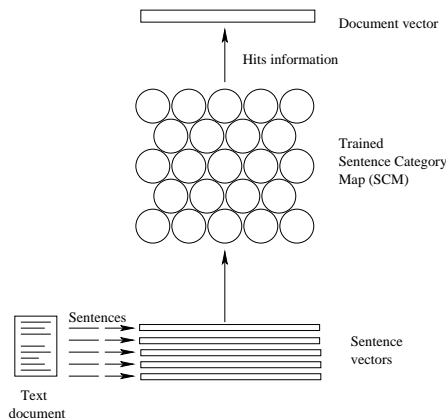


Figure 2: Extending document vectors

### III Evaluating the Goodness of Document Vectors using Information Content

It is very essential to check how well the documents are represented in terms of vectors by a particular scheme. Here we use the method proposed by [Jin *et al.*, 2001] to find the *goodness* of document vectors. The method is derived from the concepts of Latent Semantic Indexing (LSI) [Deerwester *et al.*, 1990] and information theory [Shannon, 1948]. This part describes the proposed method in brief with our experiments and results. The compelling experiments done by [Jin *et al.*, 2001] proved the intuition that *the more informative the document vectors are, the better will be the performance of IR using these vectors*. The following section provides the details.

#### 5 Mutual Information of Document Vectors

In this section we provide the necessary mathematical formulation only for showing how the information content of a document can be measured and used to find the *em* goodness of a representation. The reader is referred to [Jin *et al.*, 2001] for more details.

Let  $n$  be the number of documents in the collection. Let  $d_1, d_2, \dots, d_n$  be the document vectors in term space. Let  $M$  be the document-term matrix. Each number  $M_{ij}$  in the matrix  $M$  represents the weight

of the  $j^{th}$  word in the  $i^{th}$  document. Let  $D$  be the document-document matrix, which can be found as

$$D = MM^T \quad (3)$$

Let  $C$  be the random variable for *document content*. We define *document content* as a set of weighted *concepts* and each *concept* corresponds to an eigenvector of the document-document matrix  $D$ . Thus, the random variable  $C$  is essentially related to and can be defined in the following way: the random variable  $C$  can only take one of the values from the set of eigenvectors  $v_1, v_2, \dots, v_n$  and the eigenvalue  $\lambda_i$  indicates the importance of the eigenvector  $v_i$ . Therefore, we can assume that the probability for the random variable  $C$  to be the eigenvector  $v_i$  is proportional to the eigenvalue  $\lambda_i$ , which enables us to define the probability distribution for random variable  $C$  as the following.

$$P(C = v_i) = \frac{\lambda_i}{\sum_{j=1}^n \lambda_j}, 1 \leq i \leq n \quad (4)$$

The random variable  $D$  corresponds to the document vectors. The possible values that it can take are the set of document vectors in the document collection, *i.e.*,  $d_1, d_2, \dots, d_n$ . Since every document in the collection is equiprobable, we can assume the uniform distribution for the random variable  $D$ , that is, the probability for the random variable  $D$  to be any document vector  $d_i$  is a constant, or

$$P(D = d_i) = \frac{1}{n}, 1 \leq i \leq n \quad (5)$$

Now, the document can be viewed as a set of *concepts* and the weight for each *concept* is given by the projection of the document vector on the corresponding axis. Therefore, we can assume that the probability for a document to contain some particular *concept* is proportional to the projection of the document vector on the corresponding *concept* axis. Thus, the conditional probability  $P(C = v_i | D = d_j)$  would be proportional to the projection of document vector  $d_j$  on the *concept* axis  $v_i$ , that is:

$$P(C = v_i | D = d_j) = \frac{|d_j^T v_i|}{\sum_{k=1}^n |d_j^T v_k|} \quad (6)$$

With all these probabilities defined, we can find their respective entropies and finally, the mutual information as defined as

$$I(C, D) = H(C) - H(C|D) \quad (7)$$

The more this mutual information for a given method of vector generation, the better is that method.

## 6 Goodness using Mutual Information

We used British National Corpus (BNC) for our experiments. It is a 100 million word collection of samples from a wide range of sources, designed to represent a wide cross-section of current British English. The text corpus includes extracts from regional and national newspapers, specialist periodicals and journals for all ages and interests, academic books and popular fiction, published and unpublished letters and memorandums, school and university essays, among many other kinds of text. The detail of the documents used for our experiments is given in the following table.

| Class           | Number of docs |
|-----------------|----------------|
| Applied science | 60             |
| Arts            | 121            |
| Belief          | 69             |
| Commerce        | 92             |
| Imaginative     | 114            |
| Leisure         | 180            |
| Natural science | 49             |
| Social science  | 202            |
| World affairs   | 213            |
| Total           | 1,100          |

Table 2: BNC Documents used

| Method         | $H(C)$ | $H(C D)$  | $I(C, D)$ |
|----------------|--------|-----------|-----------|
| TFIDF          | 5.5818 | 2.8458e-6 | 5.5818    |
| Extended TFIDF | 5.6664 | 2.4029e-6 | 5.6664    |

Table 3: Mutual information for various term weighing schemes

The following observations can be made from the above results:

1. The  $H(C|D)$  value of extended TFIDF method is less than normal TFIDF method, which means that the uncertainty of guessing the document content with our method is less than the traditional TFIDF method.
2. Our method gives more mutual information showing *better* representation of the documents.

## IV Evaluation using Interclass Distance

As shown in table 2 we used 1,100 documents from nine classes of BNC. After obtaining the vector representation of these documents using TFIDF and our proposed method described in this paper, we found interclass distances for them. The results for both these schemes are summarized in figure 3. As can be seen from the results, our method is able to separate classes better than the traditional TFIDF method.



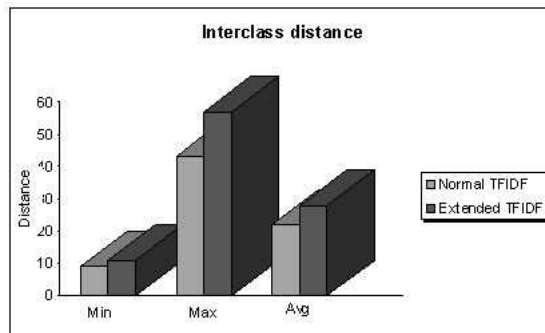


Figure 3: Interclass distance

## 7 Conclusion

In this paper we investigated the importance of using NLP for improving Information Retrieval (IR). We identified document representation as a crucial step in an IR system and found that *incorrectness* and *insufficiency* of information are the major problems with traditional methods of document representation. We proposed to use Link Grammar and some other heuristics for addressing the problem of insufficiency of the information. We showed how to extend TFIDF vectors with the help of Self-Organizing Map (SOM). In order to find how *good* these document vectors were constructed, we used their information content as a measure. Since it is shown by [Jin *et al.*, 2001] with very large scale experiments that this measure is highly correlated with the precision-recall measurements in a typical IR system, we concluded that our proposed methods do aid in improving IR. In addition to this, we also found the interclass distances for TFIDF as well as our method and further provided the support for our proposed method. We are working toward applying our method for many other corpora and IR applications.

## References

- [Agirre and Rigau, 1996] E. Agirre and G. Rigau. Word sense disambiguation using conceptual density, 1996.
- [Deerwester *et al.*, 1990] Scott C. Deerwester, Susan T. Dumais, Thomas K. Landauer, George W. Furnas, and Richard A. Harshman. Indexing by latent semantic analysis. *Journal of the American Society of Information Science*, 41(6):391–407, 1990.
- [Haykin, 1995] Simon Haykin. *Neural Networks: A Comprehensive Foundation*, chapter Self-Organizing Maps. Prentice Hall International, Inc., 1995.
- [Honkela *et al.*, 1996] T. Honkela, S. Kaski, K. Lagus, and T. Kohonen. Newsgroup exploration with WEBSOM method and browsing interface. Technical Report A32, Helsinki University of Technology, Laboratory of Computer and Information Science, Espoo, January 1996.
- [Hui, 1998] Bowen Hui. Applying NLP to IR: Why and How. Technical report, Dept. of CS, University of Toronto, April 1998.

- [Jin *et al.*, 2001] Rong Jin, Christos Faloutsos, and Alex G. Hauptmann. Meta-scoring: automatically evaluating term weighting schemes in IR without precision-recall. In *Proceedings of the 24th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 83–89. ACM Press, 2001.
- [Joachims, 1997] Thorsten Joachims. A probabilistic analysis of the Rocchio algorithm with TFIDF for text categorization. In Douglas H. Fisher, editor, *Proceedings of ICML-97, 14th International Conference on Machine Learning*, pages 143–151, Nashville, US, 1997. Morgan Kaufmann Publishers, San Francisco, US.
- [Kohonen *et al.*, 2000] T. Kohonen, S. Kaski, K. Lagus, J. Salojrvi, V. Paatero, and A. Saarela. Self organization of a massive document collection. *IEEE Transactions on Neural Networks*, 11(3):574–585, May 2000. Special Issue on Neural Networks for Data Mining and Knowledge Discovery.
- [Kohonen, 1995] T. Kohonen. *Self-Organizing Maps*. Springer, Berlin, 1995.
- [Mihalcea and Moldovan, 1998] Rada Mihalcea and Dan I. Moldovan. Word sense disambiguation based on semantic density. In Sanda Harabagiu, editor, *Use of WordNet in Natural Language Processing Systems: Proceedings of the Conference*, pages 16–22. Association for Computational Linguistics, Somerset, New Jersey, 1998.
- [Mihalcea and Moldovan, 1999] R. Mihalcea and D. Moldovan. A method for word sense disambiguation of unrestricted text. In *Proceedings of the 37th Annual Meeting of the Association for Computational Linguistics (ACL-99)*, June 1999.
- [Monz, 2000] C. Monz. Computational semantics and information retrieval. In *J. Bos and M. Kohlhase (eds.) Proceedings of the 2nd Workshop on Inference in Computational Semantics (ICoS-2), 2000*, pages 1–5, 2000.
- [Salton *et al.*, 1975] G. Salton, A. Wong, and C. S. Yang. A vector space model for automatic indexing. *Communications of the ACM*, 18:613–620, 1975.
- [Salton, 1989] Gerald Salton, editor. *Automatic Text Processing: The Transformation, Analysis and Retrieval of Information by Computer*. Addison-Wesley, 1989.
- [Sanderson, 1994] M. Sanderson. Word Sense Disambiguation and Information Retrieval. In *Proceedings of the Seventeenth Annual International ACM-SIGIR Conference on Research and Development in Information Retrieval*, pages 142–151. Springer-Verlag, 1994.
- [Shannon, 1948] Claude Shannon. A Mathematical Theory of Communication. *The Bell System Technical Journal*, pages 379–423, 623–656, July, October, 1948.
- [Sleator and Temperley, 1993] Daniel Sleator and Davy Temperley. Parsing english with a link grammar. In *Third International Workshop on Parsing Technologies*, August 1993.
- [Tokunaga and Iwayama, 1994] T. Tokunaga and M. Iwayama. Text categorization based on weighted inverse document frequency. Technical Report 94 TR0001, Department of Computer Science, Tokyo Institute of Technology, Tokyo, Japan, 1994.
- [Yang and Wilbur, 1996] Yiming Yang and John Wilbur. Using corpus statistics to remove redundant words in text categorization. *Journal of the American Society of Information Science*, 47(5), 1996.
- [Yegnanarayana, 1999] B. Yegnanarayana. *Artificial Neural Networks*, chapter 6, pages 201–232. Prentice Hall of India, 1999.