# Semantic Search Method For Textual Data

**Prof. Bhagyashree Bhoyar[1], Ankit Kumar[2], Pranjali Narkhede[3], Apeksha Nalkande[4]**

[1, 2, 3, 4] Dept of Computer Engineering

[1, 2, 3, 4] Dr.D.Y.Patil Institute of Technology, Pimpri-18

**Abstract-** *Semantic search is a methodology to improve the search accuracy by clearly understanding user search query. Normally, Most of the web search engines uses keyword based algorithm or link based algorithm that require ontology and semantic metadata to analyse user search queries . However, building a particular ontology and semantic metadata intended for large amounts of data is a costly task. In order to overcome the limitations of the conventional semantic search technique, propose a novel semantic search method that does not require ontology and semantic metadata by using the 3rd-order tensor text representation model to improve users search satisfaction. For this, it is necessary to produce the concept vector for each term that occurs in a given document, which is related to word sense disambiguation. Therefore, the semantic search that allows improving search accuracy by clearly understanding a users search intent is expected to be the core technology for next-generation search engines.*

*Keywords*- Semantic search, Text representation model, Sentence level clustering, TF-IDF

## I. INTRODUCTION

Now a days, most of online search engines are based on keyword based algorithm program or link based algorithm program. Firstly, the keyword based algorithm based program(TFIDF) searches web content by means of keywords and arranged them by date and accuracy. Secondly, the link based algorithm (page rank) arranged the web content by means of importance of the pages. It is very difficult to give the exact search result to the user queries in most of the cases of search engines. Therefore, the Linguistic search is the future technology for any search engines.

First, the information acquisition is to come up with a linguistics data from the given matter contents; the techniques like named entity recognition, and sense clarification is also used to extract linguistics data that are substantive within the given contents. Secondly, information demonstration is to make link degree without physical presence to explain number of relationship among various ideas (or knowledge nodes), and to method the linguistics search through extended ontology-based queries. metaphysics as a type of knowledgebase permits laptop machines to –exchange info with one another, and it's indispensable for developing intelligent info systems.

Finally, information utilization means that the search interface on linguistics search systems from the point of view of search users.

## II. LITERATURE REVIEW

A way for rating web content objectively and automatically, effectively measurement the human interest and a spotlight dedicated to them. Tend to compare PageRank to AN idealised random internet swimmer and tend to show a way to expeditiously work out PageRank for big numbers of pages. And, tend to show a way to apply PageRank to go looking and to user navigation [1].

Conventional information retrieval techniques With the quick expansion of web information,  are becoming insufficient for
users and often result in disappointment, because too much information can easily produce by number of keyword. Searches can be based on full-text or other content-based indexing[2].

A SemSearch, an enquiry engine, that pays special attention to the present issue by activity the complexness of linguistics search from finish users and creating it simple to use and effective. In distinction with existing semantic-based keyword search engines which usually compromise their capability of handling complicated user queries so as to beat the matter of data overhead, Semantic Search not solely overcomes the matter of data overhead however additionally supports complicated queries[3].

A design and realization of the semantic search system. The model includes three Architecture Layers of a Semantic Search System ; (they are conceptually named as) the Knowledge Acquisition, the Knowledge Representation and the Knowledge Utilization. The complexity of this algorithm is equal to O (n), it's a linear complexity, where n is number of elements and attributes in the XML document. Indeed in this algorithm a search is performed by traversing once the result Set using a loop statement[4].

Disambiguation of algorithm when an adaption of Lesk's dictionary based word sense. The source of glosses for this approach, using a standard dictionary rather. the lexical

database WordNet is employed. Provides a high ranking of semantic relations that this algorithm can utilize. [5].

A unique text house model that represents matter documents for document bunch, that contains the idea house severally of the document and term areas. The text model represented here represents documents as matrices (i.e., 2nd-order tensors), and a document corpus is depicted as a 3rd-order tensor. For this, it's necessary to provide the idea vector for every term that happens in an Page | 1 exceedingly given document, that is said to acceptation clarification. As associate external data supply for idea coefficient, tend to use the Wikipedia reference[6].

Stored documents are compared with other or with incoming patterns or documents or search requests in environment like pattern matching and document retrieval area. It appears that the best indexing (property) space is one where each entity lies as far away from the others as possible[7]

### III. RELATED WORK

**System Description:**

● I = set of input.
● O = set of output

**Process:**

**1] Query Construction**

How to construct semantic queries and how to locate their corresponding results in text representation model define as follows:

$$tj = Itj(ci) \quad \ldots(1)$$

Where,
$t_j$ = concept-by-document matrix target term.
$I_{tj}$= importance function between a target term and a concept.
$c_i$= concept space for the query term $t_j$.
C = concept space, and |C| is the total number of concepts in the tensor space model.

Let $Itj$ ($ci$) denote the importance of concept $ci$ in the concept space for the query term $tj$ is defined as follow:

$$I_{tj} (ci) = \Sigma (ci, d_k |D| ) \qquad k=1\ldots(2)$$

Where, $d_k$ represent a particular k-th document and |D| is the total number of document.

**2] Search Algorithm**

Documents are represented as 'term-by-concept' matrices and to consider the difference of document representation, need to define the similarity Function between two term-by-concept matrices of documents. By applying the *cosine* similarity function to 3 term-concept matrices, define the similarity function as follow:

$$sim(d,q) = \frac{<d,q>f}{||d2||.||q2||} \ldots (3)$$

Where, **d** represents the term-by-concept matrix of document and q represents the matrix of query**. <d,q>f** denotes the Frobenius product, which is equal to the trace of the matrices d and q.

$$<\!d,q\!>\!f = \sum_{i=1}^{n} \quad \sum_{j=1}^{m} \quad dij.qij \quad \ldots(4)$$

**3) Query Expansion with Concepts**

Associating the user's search query with the concept space, user's search intention can be expressed by clearer query. Considering the initial query, the similarity function between the document and the expanded query is defined by:

$$\text{Sim}(d_j,q)=\beta.\text{sim}(d_j,q)+(1=\beta).I_{ci}(d_j) \ldots(5)$$

Where, **sim(d_j,q)** represents the similarity between $d_j$ and original query $q$, and $I_{ci}$ $(d_j)$ represents the importance of document $d_j$ for the concept $c_i$ as equation(6).

$$\text{I}_{ci}(d_j) = \sum_{k=1}^{|T|} \quad ci(dj,tk)/N \quad \ldots(6)$$

In Equation (6), $N$ denotes the number of unique terms occurring in $dj$, which is used to normalize the importance, and $\beta$ is used to control the weight of $Ici$ $(dj)$ in the combined similarity function. When retrieving only the concept without query, $\beta$ is set to be 0.

### IV. ARCHITECTURE

The next-generation search systems can evolve thus on give a special platform and surroundings that permits users to seek out info what they require to urge a lot of simply. Propose a completely unique linguistics search technique that doesn't need ontologies and linguistics information by taking advantage of Semantically. Technique improves user's search satisfaction enriched text model.
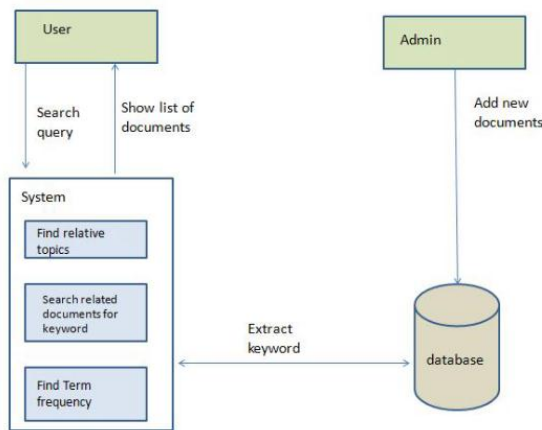
**Figure: System Architecture**

## V. EVALUATION

A Dataset have used dataset named the OSHUMED to evaluate the proposed semantic search method while constructing 3rd-order text tensor with concept-by-term-by-document. Chosen information about several categories (such as Virus Diseases, Musculoskeletal Diseases, Eye diseases, Female Genital Diseases, and Cardiovascular Diseases). Result

The goal of this experiment is to leverage the semantic search techniques and to quantitatively evaluate the degree of user satisfaction of search. It is important to remove the ambiguity of users' query for effective information retrieval, words that are commonly used are converted into the suitable query form for a specific domain through query construction method proposed in $3^{rd}$ order tensor model . Table-1 shows top three concepts for queries such as 'virus infection', 'muscle pain', 'lung cancer', 'eye disease', 'abnormal pregnancy', 'heart disorder', 'salt', and 'heart failure' in the OSHUMED dataset. For example, the terms 'Glaucoma', 'Cataract surgery', and 'Floater' are derived by 'eye disease'. A search intention would be clearer when rewriting the query with the derived concepts. Generate a text file which content summary of relevant search document.

**TABLE-1**

| Queries | Top-3 derived concepts |
|---|---|
| 'eye disease' | Glaucoma, Cataract surgery, Floater |
| 'muscle pain' | Low back pain, Anterior circulate ligament injury, Back pain |
| 'abnormal pregnancy' | Preterm birth, Abortion, Pre-eclampsia |
| 'heart disorder' | Heart failure, Health effects of tobacco, Hypertension |
| 'lung cancer' | Lung cancer, Health effects of tobacco, Acute respiratory distress syndrome |

## VI. ADVANTAGES

1. Information gets more easily according to user need.
2. Information searching requires minimum time.

## VII. DISADVANTAGE

Huge database leads more time consumption to get the information.

## VIII. APPLICATION

1. Web search engine
2. Searching the content from any databases

## IX. CONCLUSION

The next-generation search systems can evolve thus on offer a special platform and atmosphere that permits users to search out data what they require to urge a lot of simply. During this context, linguistics search technology is going to be the core main mechanism of the next-generation of search engines as a result of it's the simplest thanks to minimizing the psychological feature effort of users and to satisfy their data wants. In developing semantic search services, it's an great challenge to create a helpful results without physical presence. During this paper, we tend to plan a brand new semantic search technique with the Wikipedia-based 3rd-order tensor text model, that doesn't need developing an exact metaphysics. Through intensive experiments mistreatment the OSHUMED document assortment and SCOPUS library information, tend to evidence that the planned ways improve users' search satisfaction fairly.

## REFERENCES

[1] L. Page, S. Brin, R. Motwani, and T. Winograd, Page, Lawrence, et al. The PageRank citation ranking: bringing order to the Web, 1999.

[2] F. Wissbrock, Information Need Assessment in Information Retrieval; Beyond Lists and Queries, Proceedings of the 27th German Conference on Artificial Intelligence, 2004.

[3] Y. Lei, V. Uren, Y. Kanza, and Y. Sagiv, Semsearch: A search engine for the semantic web, Managing Knowledge in a World of Networks, Springer Berlin Heidelberg, pp.238-245, 2006.

[4] D. I. Hana, H. I. Kwonb, and H. J. Chong, A Study on the Conceptual Modeling and Implementation of a Semantic Search System, Journal of Intelligence and Information Systems, Vol. 14, No. 1, pp.67-84, 2012.

[5] S. Banerjee and T. Pedersen, An adapted lesk algorithm for word sense disambiguation using word net, 2015.

[6] H. J. Kim, K. J. Hong, and J. Y. Chang, Semantically enriching text representation model for document clustering, Proceedings of the 30th Annual ACM Symposium on Applied Computing, pp. 922-925. 2015.

[7] G. Salton, A. Wong, and CS. Yang, A vector space model for automatic indexing,