

## DEEP WEB SEARCHING (DWS)

Chandru V<sup>1</sup>, Dinesh V<sup>2</sup>, Jeffrin J<sup>3</sup>, Pradeep K.R<sup>4</sup>, Sharmila D<sup>5</sup>

<sup>1,2,3,4</sup>B.Tech Information Technology, Dr.NGP Institute of Technology, Coimbatore.

<sup>5</sup>HOD, Dept. of Information Technology, Dr.NGP Institute of Technology, Coimbatore, Tamilnadu.

\*\*\*

**Abstract** - As deep web grows at a very fast pace, there has been increased interest in techniques that help efficiently locate deep-web interfaces. An exploratory search may be driven by a user's curiosity or desire for specific information. When users investigate unfamiliar fields, they may want to learn more about a particular subject area to increase their knowledge rather than solve a specific problem. A matching query style has significant limitations. Search results are satisfactory only when users give the right search words. To achieve more accurate results for an exploratory search, Smart Crawler ranks websites to prioritize highly relevant ones for a given topic.

**Key Words:** Document clustering system, Fuzzy-logic, Self-Organized Mapping (SOM) algorithm, Ontology, Text mining, Wordnet.

### 1. INTRODUCTION

Individuals with disabilities understands the actual content of the web page in a more efficient manner. Text clustering is mainly used for a document clustering system which clusters the set of documents based on the user typed key term. First the system preprocesses the set of documents and the user given terms. The feature evaluation is used to reduce the dimensionality of high-dimensional text vector. Proposed a fuzzy-logic based model as a decision tool for results selection. The new proposals in each discipline are clustered using a SOM algorithm.

#### 1.1 OBJECTIVE

Text clustering is mainly used for a document clustering system which clusters the set of documents based on the user typed key term. Firstly, the system preprocesses the set of documents and the user given terms. The feature evaluation is used to reduce the dimensionality of high-dimensional text vector. Proposed a fuzzy-logic based model as a decision tool for results selection. The new proposals in each discipline are clustered using a SOM algorithm.

#### 1.2 WORDNET

WorldNet® is a large lexical database of English. Nouns, verbs, adjectives and adverbs are grouped into sets of cognitive synonyms (synonym sets), each expressing a distinct concept. Synonym sets are interlinked by means of conceptual-semantic and lexical relations. The main relation

among words in WordNet is synonym, as between the words shut and close or car and automobile. Synonym words that denote the same concept and are interchangeable in many contexts are grouped into unordered sets (synonym sets).

### 2. SYSTEM OVERVIEW

To cluster the text documents over the web page based on the user typed key term. To enhance deep web search (ontology) and overcome grouping of unrelated documents into the same cluster. Aims to help web users locate the best search tools for their search needs, resulting in faster and more accurate search results. Present work assumes that all user local instance repositories have content-based descriptors referring to the subjects, however, a large volume of documents existing on the web may not have such content-based descriptors. For this problem strategies like ontology mapping and text classification/clustering were suggested. These strategies will be investigated in future work to solve this problem.

### 3. SOFTWARE OVERVIEW

**Java is a platform Independent.** Java is a high-level programming language Introduced by Sun Microsystems in June 1995 Java is becoming a standard for Internet Applications. It provides for interactive processing and for the use of graphics and animation on the Internet. Since the Internet consists of different types of computers and operating systems, a common language was needed to enable computers to run programs that run on multiple platforms. Java is an object-oriented language built upon C and C++.It derives its syntax from C and its object-oriented features are influenced by C++. Java can be used to create applications and applets. An application is a program that runs on the user's computer, under its operating system. An applet is a small window-based program that runs on HTML page using Java enabled We browser like Internet Explorer, Netscape Navigator, Hot Java or an applet view

### 4. SYSTEM IMPLEMENTATION

Implementation is the most crucial stage in achieving a successful system and giving the user's confidence that the new system is workable and effective. It may be implementation of a modified application to replace

an existing one. This type of conversation is relatively easy to handle, provide there are no major changes in the system.

Each program is tested individually at the time of development using the data and has verified that this program linked together in the way specified in the programs specification, the computer system and its environment is tested to the satisfaction of the user. The system that has been developed is accepted and proved to be satisfactory for the user. And so the system is going to be implemented very soon. A simple operating procedure is included so that the user can understand the different functions clearly and quickly.

Initially as a first step the executable form of the application is to be created and loaded in the common server machine which is accessible to the entire user and the server is to be connected to a network. The final stage is to document the entire system which provides components and the operating procedures of the system. Implementation is the stage of the project when the theoretical design is turned out into a working system. Thus, it can be considered to be the most critical stage in achieving a successful new system and in giving the user, confidence that the new system will work and be effective.

The implementation stage involves careful planning, investigation of the existing system and its constraints on implementation, designing of methods to achieve changeover and evaluation of changeover methods. Implementation is the process of converting a new system design into operation. It is the phase that focuses on user training, site preparation and file conversion for installing a candidate system. The important factor that should be considered here is that the conversion should not disrupt the functioning of the organization.

### 5. SYSTEM ARCHITECTURE

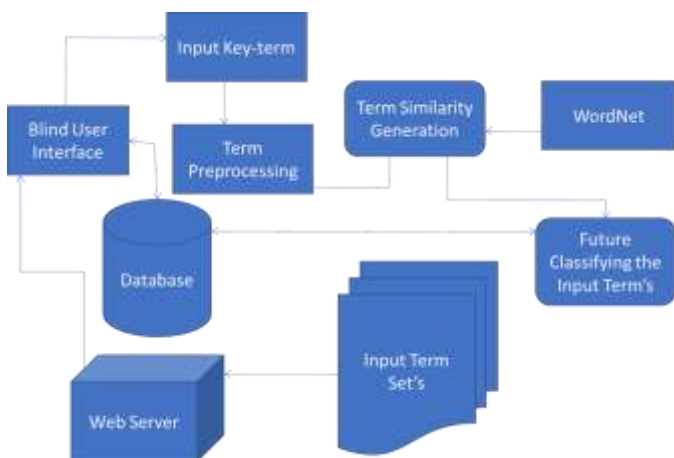


Figure 5.1 System Implementation

Implementation is the stage of the project when the theoretical design is turned out into a working system. Thus,

it can be the considered to be the most critical in achieving a successful new system and in giving the user, confidence that the new system will work and be effective. The implementation stage involves careful planning, investigation of the existing system and its constraints on implementation, designing of methods to achieve changeover and evaluation of changeover methods.

### 6. DATA FLOW DIAGRAM



Figure 6.1 DF Diagram

### 7. STOP WORD REMOVAL:

Stop words are words which are filtered out prior to, or after, processing of natural language data (text). It is controlled by human input and not automated. These are some of the most common, short function words, such as the, is, at, which and on.

### 8. Text Analysis:

The Artificial-Intelligence literature contains many definitions of ontology (Word net).

- It includes machine-interpretable definitions of basic concepts in the domain and relations among them.
- The featured results produced by the sentence-based, document-based, corpus-based, and the combined approach concept analysis have higher quality than those produced by a single-term analysis similarity.

### 9. MULTI-TERM SEARCH

Get the multi-term input from the user and it will search the keyword one by one and get the relevant content from the web servers. Our system get the search result deeply from the search engines and its search the terms randomly till last key term in that multi-term list.

## 10. CONCLUSION

This paper has presented an OTMM for grouping of research proposals. Research ontology is constructed to categorize the concept terms in different discipline areas and to form relationships among them. It facilitates text-mining and optimization techniques to cluster research proposals based on their similarities and then to balance them according to the applicants' characteristics. The experimental results at the NSFC showed that the proposed method improved the similarity in proposal groups, as well as took into consideration the applicants' characteristics (e.g., distributing proposals equally according to the applicants' affiliations).

## REFERENCES

- [1] T. Deng, L. Zhao, H. Wang, Q. Liu, and L. Feng. Refinder: a context-based information re-finding system. *IEEE TKDE*, 25(9):2119–2132, 2013.
- [2] J. Hailpern, N. Jitkoff, A. Warr, K. Karahalios, R. Sesek, and N. Shkrob. Youpivot: improving recall with contextual search. In *CHI*, pages 1521–1530, 2011.
- [3] J. Teevan, E. Adar, R. Jones, and M. Potts. Information retrieval: repeat queries in yahoo's logs. In *SIGIR*, pages 151–158, 2007.