

Plagiarism Scan Report	
Summary	
Report Generated Date	03 Nov, 2017
Plagiarism Status	100% Unique
Total Words	543
Total Characters	3586
Any Ignore Url Used	

Content Checked For Plagiarism:

Abstract— In Semantic Web, Text Summarization, Search Engine Optimization and many more such technologies, document classification is a key aspect in getting the results. Extracting domain keywords from documents helps to optimize the task of document classification involved in Information Retrieval. The existing state of the art techniques extensively depend on keyword extraction based on term document frequency. Also, the existing technologies rank words based on the title of the document, which, in some cases, is imprecise because sometimes the title of the document doesn't have words relevant to the context in mind. To overcome such problems, we propose the idea of ontology based keyword extraction for increasing the accuracy of document classification and in turn its applications. The objective of this paper is to extract domain specific keywords from the given text document with the help of Domain Dictionary created using Ontology. This approach can be further extended towards revitalizing text summarization techniques.

Keywords— ontology; domain dictionary; keyword extraction; information retrieval; term document frequency

1. INTRODUCTION

Document classification is required for many applications such as Search Engine Optimization, Semantic Web, Metadata Tagging, etc. Existing applications or APIs for doing so for a given domain could not be found. Furthermore, the ambiguity in existing applications for applications of documents classification is high, which results in incorrect data and thus, incorrect results. Moreover, Term Document Classification is used extensively in applications of Information Retrieval. In some cases, this results in data which is not in the context. This method needs to be refined.

Thus, there is a need for coming up with a new methodology or implementation of the same, and we are proposing it through Ontology. Ontology will help us improve the accuracy of these applications to a great extent by giving

results based on the context of the query.

Ontology is a way of knowledge representation in which all the concepts related to a particular object are explored and relations are established between these concepts. Ontology should be such that it is machine readable which enables it to comprehend it and it should be expressed in some manner.

Ontology is a formal description of concepts, properties of these concepts which are basically the features of the concepts and restriction imposed on these properties. To develop an ontology we need to first define all the classes (concepts) in the domain under consideration, then identifying classes as superclass or subclass and the relationship between these classes. Afterwards, we need to define the slots (properties) and the restriction imposed on them that is permissible values or each slot. All this constitutes the knowledge base.

When two parties are communicating, both should have already established the context of the communication. For e.g., if party 'A' talks about 'Jaguar' with party 'B', then B must know whether A is talking about Jaguar the animal or car or operating system. Ontology establishes this context between both the users.

In this paper we propose a system where the user provides the document(text) and the domain of concern, the system extract keywords with help of pre-created ontology and gives all keywords along with important parameters such as frequency of word, how strong is the association between the word and domain, etc. which are discussed in greater detail in the paper.

Report generated by smallsecrets.com

Plagiarism Scan Report	
Summary	
Report Generated Date	03 Nov, 2017
Plagiarism Status	100% Unique
Total Words	536
Total Characters	3643
Any Ignore Url Used	

Content Checked For Plagiarism:

LITERATURE SURVEY

International scholars have studied this field for many years now and explored various methods. Paper [1] proposes a keyword dictionary server that provides keyword expansion using domain specific ontologies. This has been achieved using the functional metadata of services like service name, category, provider and description. Paper [2] propose the skeleton of a semantic search engine that follows automatic

query expansion. For all the terms, SPARQL query is built and then it is fired on the knowledge base that finds appropriate RDF triples in knowledge Base. Web documents relevant to the requested concepts and individuals specified in these triples are then retrieved and ranked according to their relevance to the user's query and then are sent to the user. Paper [3] uses WordNet as a dictionary for finding synonyms of user's query. This paper explores a technique called ontological indexing which is based on calculating the context of the words using ontology. In paper [4], ontology is created by domain experts and is supplied to the system. Here 2 algorithms are proposed for extraction: "semantic information extracting algorithm" and "semantic information re-recognition algorithm". Text information is extracted using ontology and the 2 proposed algorithms.

Paper [5] talks about using Ontology Based Information Extractors(OBIE) for text grading. They argue that the combination of information extractors that perform different functions can provide a better understanding of a graded text, and the combination of information extractors that have different implementations can improve the performance of the extraction process. Paper [6] enables the ontology to find related recent knowledge in the domain from communities, by exploiting their underlying knowledge as keywords. It extracts instances and statements from the documents using the ontology-based and pattern-based information extraction technique. A confidence value is used in order to maintain the stability of the ontology. Finally, the proposed system enriches

the ontology with the new extracted instances and statements and validates the knowledge inside the ontology. Paper [8] reviewed the related concepts and methods of ontology construction and extension, proposed an automatic ontology extension method based on supervised learning and text clustering. Paper [9] proposes an approach to extract ontology directly from RDB in the form of OWL/RDF triples, to ensure its availability at semantic web. Then it automatically constructs an OWL ontology from RDB schema using direct mapping rules. Later, rewriting SPARQL query from SQL by translating SQL relational algebra into an equivalent SPARQL. In paper [7], authors have developed a framework for comparing 11 ontology learning systems. After analyzing these methodologies, 3 are selected and then applied to finance domain. Out of the 11 ontology learning systems, we have selected 5 systems : OntoLearn [10] uses text mining and statistical techniques to learn concepts and build taxonomic relations; Text2Onto [11] uses Probabilistic Ontology Model and involves statistical, linguistic techniques to create ontology; CRCTOL [12] algorithm is an statistical algorithm, which extracts concepts and relations; OntoGain [14] is an unsupervised algorithm which uses linguistic tools to preprocess text and extract concepts; HCHIRISM [13] is also an unsupervised algorithm which recursively analyze a large number of web sites in order to find important concepts for a domain by introducing an initial keyword. A comparative study of all these methods is put in Appendix.

Plagiarism Scan Report	
Summary	
Report Generated Date	03 Nov, 2017
Plagiarism Status	100% Unique
Total Words	467
Total Characters	2935
Any Ignore Url Used	

Content Checked For Plagiarism:

PROPOSED MODEL

The two main components of the proposed model are:

- i) Domain Ontology Creation
- ii) Parameter application on text

The Domain Ontology creation module generates an Ontology based on the [13] HCHIRSIM model.

The second module shortlists the important keywords from the text using a traditional feature-based selection approach.

Figure 1. demonstrates the system design incorporating the two modules to accomplish the goal of Domain specific keywords extraction.

The two main inputs to the system are i) User's document from which words are to be extracted and ii) The desired domain. Of these the domain is given as input to the Ontology Creation Module.

The proposed model works as follows:

1. The Domain Ontology which is used is created.
2. The user's document is pre-processed and cleaned.

The two pre-processing steps are Cleaning and Lemmatization

- a. Cleaning involves the removal of special characters and stopwords from the text.
- b. Lemmatization is used to derive the root form of the word from its inflection version.
3. The words from this preprocessed text are then mapped to concepts and instances in the Ontology created in step 1. Mapping involves looking up the existence of the words in preprocessed text in the Ontology.
4. Parallely the keywords from the original document are shortlisted based on the following parameters:
 - i. Frequency of the word in the English language.
 - ii. Frequency of occurrence in the text.

- iii. Position weightage.
 - iv. Part of Speech of the word.
 - v. Number of times the word is used as either a Subject, Object or a Predicate.
 - vi. Distance between current and previous occurrence.
5. The Keywords obtained from the Ontology Mapping and Parameter Application are then ranked.
 6. Finally the highest ranked keywords are output as the Domain specific keywords.

IV. CONCLUSION

In this paper, we have explained the need for Ontology in identifying keywords in a given document and have understood how it is useful in various disciplines. We highlighted that the existing systems use NLP and often extract keywords based on title. The delta change in the proposed system is to extract keywords based on domain. Our proposed system has two main modules which are Ontology creation module and Keyword extraction module. By employing Ontology we can overcome the ambiguity that the existing system suffers from.

Major application of proposed system is Classification of documents, which further can be applied to Search engine optimization, Semantic web etc.

ACKNOWLEDGMENT

We are thankful to our college Vivekanand Education Society's Institute of Technology for considering our project and extending help at all stages needed during our work of collecting information regarding our project. We are deeply indebted to Head of the Computer Department Dr.(Mrs.) Nupur Giri and our Principal Dr. (Mrs.) J.M. Nair, for giving us this valuable opportunity to do this project.