

Creation and Enhancement of Fuzzy Ontology Structures for Query Answering from Text Documents – A Text Mining Paradigm

by

Muhammad Abulaish

Department of Mathematics

Submitted

in fulfillment of the requirements of the degree of Doctor of Philosophy

to the



Indian Institute of Technology Delhi

April 2006

LIBRARY
No. T.H. 3344

I.I.T. DELHI
LIBRARY
PROCESSED

To my parents, who taught me how to tackle life's many impediments.


*To my wife, sons – Yusuf and Haris, and daughter – Aisha who make the
journey worthwhile.*

Certificate

This is to certify that the thesis entitled "**Creation and Enhancement of Fuzzy Ontology Structures for Query Answering from Text Documents – a Text Mining Paradigm**" submitted by Mr. Muhammad Abulaish to the Indian Institute of Technology Delhi, for the award of the degree of Doctor of Philosophy, is a record of the original bona fide research work carried out by him under my supervision and guidance. The thesis has reached the standards fulfilling the requirements of the regulations relating to the degree.

The results contained in this thesis have not been submitted in part or full to any other university or institute for the award of any degree or diploma.

New Delhi
April 2006


Dr. Lipika Dey
Supervisor

Acknowledgements

I am very grateful to my supervisor, Dr. Lipika Dey, for having accepted me as a student and introducing me to such an interesting field. With her sagacious suggestions, understanding, patience, enthusiasm, responsibility, vivid discussions, insight, sensibility, observation, rigor, and above all with her appreciating nature, she helped me to work in a right perspective. I cannot thank her enough for the time and effort she spent in disciplining my style of working, introducing me to new problems and helping me put things in right perspective, professionally and personally. I am deeply indebted to her for her intense involvement, effort, advice, encouragement, patience and help. Her deep insight into problems, imaginative ideas, zest for work, foresight, rigour and tremendous capacity for disciplined work, have been a source of inspiration and motivation for me. Working with her has been a day-to-day learning and invigorating experience. In spite of being extremely busy, she could always find time for my work; means so much to me. Without her planning, immaculate work scheduling, active interest and positive thinking, my thesis would not have been through.

I take this opportunity to thank Prof. B. R. Handa who selected me and gave me a chance to get enrolled in the department for the Ph. D. program. I express my regards to Prof. B. Chandra, DRC Chairperson and Head of the Department, for her support. Special thanks to my SRC Members: Prof. Suresh Chandra, Dr. Nesar Ahmad and Dr Tapan Roy Chowdhury for spending their valuable time during the discussions over seminars. I am grateful to my pre-Ph.D. course instructors: Prof. B. Chandra, Dr. Lipika Dey, Dr. Wagesh Shukla, Dr. Niladri Chatterjee and Dr. S. Chandra Shekhar Rao for their valuable suggestions as well. Further, I extend my thanks to Dr. R. K. Sharma, Dr. S. Dharmaraja and other faculty members of the department for their encouragement. I am also thankful to Mrs. Rashmi

Wadhwa and other non-teaching staff members of the department for all their help and good wishes. I am also thankful to I.I.T. Delhi authorities for providing me the necessary facilities.

It is a pleasure for me to thank Dr. Shailendra singh and Dr. K. V. Krishna for their suggestions to get acquainted with the department. I am also thankful to Mr. Sachin Kumar and all other research scholars of the department for their company.

I am very grateful to the Department of Science and Technology, Ministry of Science and Technology, who sponsored me to attend the CIMPA-UNESCO-INDIA school on "Soft Computing Approach to Pattern Recognition and Image Processing", held at Indian Statistical Institute, Kolkata during December 2-13, 2002 which was the foundation on which I started my research work.

I would also like to thank Dr. Siraj Hussain (IAS), ex Vice-Chancellor of Jamia Hamdard (Deemed to be University) for his encouragement and permission to get enrolled for the Ph.D. program.

I am also thankful to my elders and colleagues, especially Prof. Sharfuddin Ahmad, Prof. M. R. Khan, Dr. Naseem Ahmad, Dr. Shehzad Hasan and other staff members of my working department for their kind support without which it would not be possible for me to finish this work on time. I owe special thanks to my friend Mr. Zahiruddin for his kind support.

To my family, I owe so much more than words can express, for what **they have done for me**. It was only because of them that I could do a Ph. D. Their constant encouragement, love, support and advice have been so important. For being there for me at all times and for taking over all the problems so that I could concentrate single-mindedly on my research.

Finally, I would like to thank almighty God who helped me at different critical stages during this work and shown me the right way when I was in trouble.

New Delhi

Md. Abulaish

Muhammad Abulaish

Abstract

Over the last decade the World Wide Web (WWW) has virtually become the main repository for data sharing and information gathering. Most of the data on the Web are embedded within text documents which are either unstructured or semi-structured in nature. Thus even though there is no scarcity of information on the Web, locating, extracting and analyzing required information from this vast unstructured collection is a complex and challenging task. Most of the problems are obviously due to the incapability of the computer in comprehending natural language texts with all its nuances.

Word based searches for relevant information from texts retrieve a huge collection and burden the user with information overload. Ontology based text information retrieval can perform concept-based search and extract only relevant portions of text containing concepts that are present in the query or those that are semantically linked to query concepts. While these systems have better precision of retrieval than general-purpose search engines, there are some crucial problems which have to be tackled while designing such systems. Considering that ontologies are created and maintained by domain experts, creation and maintenance of ontologies are expensive tasks. Though ontology provides a structured framework to store domain knowledge, all ontological concepts cannot be unambiguously described using crisp property descriptors. Besides, the ontological descriptors may not exactly match text descriptions or the user given descriptors in query. Appropriate reasoning mechanisms have to be designed in order to apply ontologies for text information processing.

In this thesis we have proposed the use of a text mining framework which can integrate the problems of ontology enhancement with text information retrieval. Thus the system provides

dual advantage of content-based retrieval and ontology learning from text documents. In order to achieve flexibility in concept definition we propose to enhance the traditional ontology structures into fuzzy ontology structures. These structures can accommodate imprecise concept descriptors that are prevalent in text documents. The ontology-based text mining framework mines precise and imprecise concept descriptions from text documents and enhances the fuzzy ontology structure with them. The mined knowledge is stored in a structured knowledge base which design is derived from the underlying fuzzy ontology structure representing the domain. We have shown that the fuzzy ontology structure can provide a viable solution to resolve inconsistencies over concept descriptions encountered in multiple overlapping ontologies. The proposed framework is also enriched with a fuzzy reasoning mechanism that is used to process user queries over a curated knowledge base.

The text-mining framework is enhanced to mine inter-concept relations besides structural semantic relations from text. The query processing principle is then appropriately extended to answer relation-based complex user queries over a curated database. This has been tested over biomedical document repositories, in which relations play an important role in focused information retrieval. Since mining can yield a wide array of relations, these relations are subjected to a significance analysis to identify the key relations of a domain. The significant relations are used for enhancing the underlying concept ontology into a relational ontology.

Table of Contents

| | |
|---|-----------|
| Certificate | i |
| Acknowledgement | iii |
| Abstract | v |
| List of Figures | xiii |
| List of Tables | xvii |
| 1 Introduction | 1 |
| 1.1 Text Information Retrieval | 1 |
| 1.2 Approaches to Overcome the Problem of Information Overload | 3 |
| 1.2.1 Ontology-based Text Processing | 5 |
| 1.3 Shortcoming of Existing Ontology Designs and Applications | 6 |
| 1.4 Scope of the Proposed Work | 11 |
| 1.5 Overview of the Thesis | 13 |
| 2 From Search Engines to the Semantic Web: Review of Current Trends in Text Information Processing | 17 |
| 2.1 Introduction | 17 |
| 2.2 Core Concepts in Text information Processing | 19 |
| 2.3 Evolution of Text Information Retrieval | 22 |
| 2.3.1 Search Engines | 23 |
| 2.3.2 Improving Retrieval Performance | 25 |
| 2.4 Content-based Text Retrieval | 26 |
| 2.4.1 Information Extraction (IE) | 27 |
| 2.4.2 Curation | 27 |
| 2.4.3 Query Answering | 28 |
| 2.4.4 Question Answering | 28 |
| 2.5 Semantic Web | 30 |

| | | |
|----------|---|-----------|
| 2.6 | Ontology | 32 |
| 2.6.1 | Ontology Representation | 33 |
| 2.6.1.1 | eXtensible Markup Language (XML) | 33 |
| 2.6.1.2 | Resource Description Framework (RDF) | 34 |
| 2.6.1.3 | Resource Description Framework Schema (RDFS) | 35 |
| 2.6.1.4 | DAML ONTology (DAML-ONT) Language | 36 |
| 2.6.1.5 | Ontology Inference Language (OIL) | 36 |
| 2.6.1.6 | DAML+OIL | 37 |
| 2.6.1.7 | Web Ontology Language (OWL) | 38 |
| 2.6.2 | Ontology Engineering Tools | 40 |
| 2.6.3 | Ontology Reasoning | 42 |
| 2.7 | Popular Ontologies and their Applications | 43 |
| 2.8 | Ontology-based Text Information Processing | 47 |
| 2.9 | Ontology Learning and Enhancement | 49 |
| 2.9.1 | Integrating Ontology Learning with Text Information Extraction | 50 |
| 2.9.2 | Introducing New Relations into an Ontology | 52 |
| 2.10 | Text Mining for Ontology Enhancement | 54 |
| 2.11 | Fuzzy Ontology Structure | 56 |
| 2.12 | Integrating Multiple Domain Ontologies | 57 |
| 2.13 | Conclusion | 63 |
| 3 | A Fuzzy Ontology Structure for Imprecise Concept Descriptions | 65 |
| 3.1 | Introduction | 65 |
| 3.2 | Ontology for Structured Representation of Domain Knowledge..... | 68 |
| 3.3 | The Problem of Handling Imprecise Concept Descriptions | 72 |
| 3.4 | Proposed Fuzzy Ontology Structure for Including Imprecise Concept Descriptions | 77 |
| 3.4.1 | Encoding Domain Knowledge using Fuzzy Ontology Structures | 84 |
| 3.4.2 | Creation of Fuzzy Ontology Structure | 92 |
| 3.5 | Handling Inconsistent Ontology Concept Descriptions | 98 |
| 3.6 | Fuzzy Ontology Structures for Representing Inconsistent Concept Descriptions | 100 |

| | | |
|----------|--|------------|
| 3.6.1 | Computing Concept Consistency | 104 |
| 3.7 | Conclusion | 107 |
| 4 | Database Curation and Ontology Enhancement through Ontology-based Text Mining | 109 |
| 4.1 | Introduction | 109 |
| 4.2 | The Text Mining Framework | 112 |
| 4.3 | Proposed Ontology-based Text Mining System | 114 |
| 4.3.1 | Ontology Parser | 117 |
| 4.3.2 | Document Processor | 121 |
| 4.3.3 | Knowledge Distiller | 125 |
| 4.4 | Experimental Results for Text Mining and Ontology Enhancement | 131 |
| 4.5 | Performance Analysis of Proposed Text Mining System | 144 |
| 4.6 | Conclusion | 148 |
| 5 | Imprecise Query Answering from Curated Databases | 149 |
| 5.1 | Introduction | 149 |
| 5.2 | Imprecise Query Processing | 151 |
| 5.3 | Imprecise Query Processing System | 153 |
| 5.3.1 | Query Interface Design and Query Formulation | 155 |
| 5.3.2 | Query Parser and SQL Generator | 157 |
| 5.3.3 | Fuzzy Query Processor..... | 158 |
| 5.4 | Imprecise Query Processing Examples | 162 |
| 5.5 | Evaluation of the Query Answering Process | 172 |
| 5.6 | Conclusion | 173 |
| 6 | Mining Relations from Ontologically Tagged Texts – Application to Biomedical Domain | 175 |
| 6.1 | Introduction | 175 |
| 6.2 | Ontology-based Text Processing for Biomedical Domain | 179 |
| 6.3 | GENIA Ontology for Molecular Biology | 182 |
| 6.4 | GENIA Corpus | 184 |
| 6.5 | Generic Biological Relation and Related Concepts | 185 |

| | | |
|----------|---|------------|
| 6.6 | Generic Biological Relation Mining System | 188 |
| 6.6.1 | Document Processor | 191 |
| 6.6.2 | Biological Relation Extractor | 197 |
| 6.6.2.1 | Extraction of Information Components | 197 |
| 6.6.2.2 | Identifying Feasible Generic Biological Relations | 200 |
| 6.6.3 | Feasible Generic Biological Relations Identified from the GENIA Corpus | 203 |
| 6.6.4 | Generating Feasible Fuzzy Biological Relations | 205 |
| 6.7 | System Performance Evaluation | 208 |
| 6.8 | Uniqueness of the Proposed System over Earlier Systems | 212 |
| 6.9 | Conclusion | 214 |
| 7 | Database Curation and Relation-based Query Answering for Biomedical Domain | 217 |
| 7.1 | Introduction | 217 |
| 7.2 | Database Curation | 218 |
| 7.3 | Maintaining Knowledge Base for Query Answering | 223 |
| 7.3.1 | Knowledge Base Indexing | 223 |
| 7.4 | Relation-based Query Answering for Biomedical Domain | 226 |
| 7.4.1 | Example Queries and their Results | 229 |
| 7.5 | Performance Evaluation of the Relation-based Query Answering System | 237 |
| 7.6 | Conclusion | 239 |
| 8 | Enhancing Biological Ontologies to Fuzzy Relational Ontology Structures | 241 |
| 8.1 | Introduction | 241 |
| 8.2 | Characterizing Generic Biological Relations | 243 |
| 8.3 | Generating a Concept-pair Tree to Define Biological relations | 246 |
| 8.4 | Mapping the Relation Instances over a Concept-pair Tree | 251 |
| 8.5 | Computing Information Loss at Concept-pair Nodes | 254 |
| 8.6 | Relations Strength Calculation and their Fuzzification | 257 |
| 8.7 | Fuzzy Relational Ontology Model | 263 |
| 8.8 | Enhancing GENIA to a Fuzzy Relational Ontology | 264 |

| | |
|---|------------|
| 8.9 Conclusion | 265 |
| 9 Conclusion | 267 |
| 9.1 Contributions of the Thesis | 267 |
| 9.2 Extensions to the Proposed System..... | 271 |
| 9.3 The Road Ahead - Future Directions in Ontology-based Text Information Retrieval | 272 |
| 9.4 Lessons Learnt and Pitfalls to be Avoided | 273 |
| Bibliography | 275 |
| Appendices | 291 |
| A1 OWL Codes for Fuzzy Wine Ontology Structure | 291 |
| A2 Stop Words used by PubMed Search Engine | 295 |
| A3 List of Feasible Biological Relation and their Morphological Variants Mined from GENIA Corpus | 296 |
| A4 Biological Relation Triplets and their Strengths | 298 |
| Brief Bio-data of the Author | 307 |